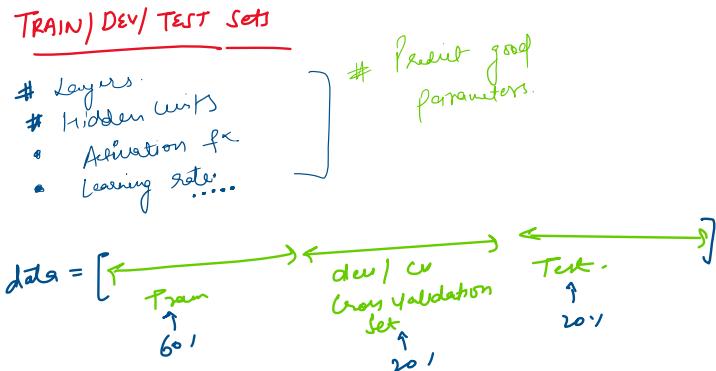


## Setting up your Machine Learning Application.

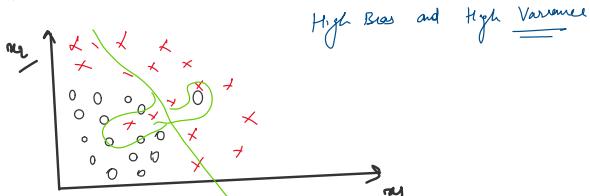


# Not having a test set is okay and only have Dev set.

### Bias / Variance

Train set Error  $\rightarrow 17$ .  
Dev set Error  $\rightarrow 117$ .

Underfitting - High Bias



### Basic Recipe for Machine Learning

- ① High Bias ??  $\rightarrow$  \* Bigger N/w Layer  
(Train set performance)  
• Layer optimization Algo's
- ② High Variance ?  $\rightarrow$  \* More data  
(Dev set performance)  
• Regularization  
• neural N/w Architectures

Done!

### bias / Variance Tradeoff

bias  $\uparrow$  Variance  $\downarrow$

## Regularizing your Neural N/w

$$\|w^u\|^2 = \sum_{i=1}^m \sum_{j=1}^{n(u)} (w_{ij}^{(u)})^2$$

### # Regularization

$$\text{Logistic Regression} - \min_{(w,b)} J(w,b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2} \|w\|^2$$

L<sub>2</sub> Regularization  $\rightarrow w \in \mathbb{R}^{n \times 1}, b \in \mathbb{R}$

$$\|w\|^2 = \sum_{j=1}^n w_j^2 = w^T w$$

$$\|w\|_1 = \sum_{j=1}^n |w_j| \quad (\lambda = \text{Regularization parameter})$$

L<sub>2</sub> Regularization

$$\frac{1}{m} \sum_{i=1}^m \|w\|_2^2 = \frac{\lambda}{m} \|w\|_F^2$$

( $\lambda$  = Regularization parameter)  
[Hyperparameter]

Neural N/w

$$J(w^{(0)}, b^{(0)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{(l+1)}} (w^{(l)}_{ij})^2$$

} "Frobenius norm"  $\| - \|_F$

$$w \in (\mathbb{R}^{n^{(l)}}, \mathbb{R}^{n^{(l+1)}})$$

$$dw = (\text{Calculated from Backprop}) + \frac{\lambda}{m} w$$

$$\frac{\partial J}{\partial w^{(l)}} = dw^{(l)}$$

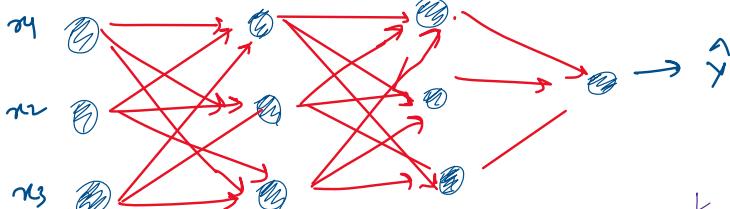
$$w^{(l)} = w^{(l)} - \alpha dw^{(l)}$$

# L<sub>2</sub> regularization is also called "weight decay".

$$w^{(l)} := w^{(l)} - \alpha \left[ (\text{Backprop}) + \frac{\lambda}{m} w^{(l)} \right]$$

$$= w^{(l)} - \frac{\alpha \lambda}{m} w^{(l)} - \alpha (\text{Backprop})$$

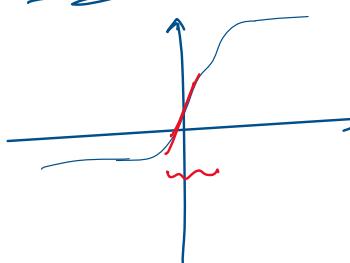
Why Regularization prevents Overfitting?



$$J(w^{(0)}, b^{(0)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

- ① (If  $\lambda \gg$  big,  $w \approx 0$ )
- ② high bias  $\rightarrow$  high variance

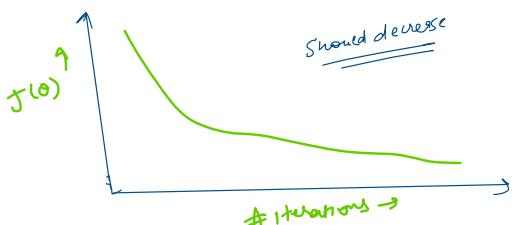
tanh(f(x))



$g(x) = \tanh(x)$   
 $w^{(l)} \downarrow, \lambda \uparrow$   
 $f(x)$  will be linear  
in red range

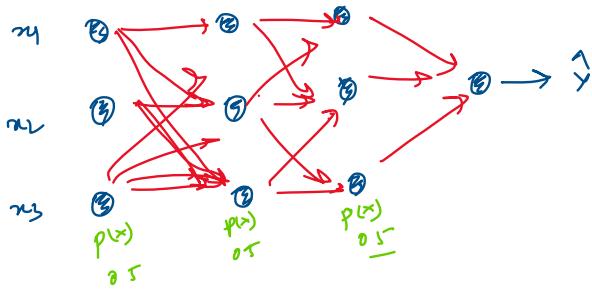
# As if every layer is linear  $N/w$

$$J(\cdot) = \sum L(\hat{y}, y) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$



## Dropout    REGULARIZATION

④ Set some probability to dropout a node in now,



## Implementing Dropout (Inverted dropout)

# illustrate in layer  $L=3$ :

~~# Reshape~~  
~~d3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep\_prob.~~  
~~keep prob = 0.8 (say)~~

$$a_3 = \text{np multiply } (a_1, d_3) \quad (a_3 \neq d_3)$$

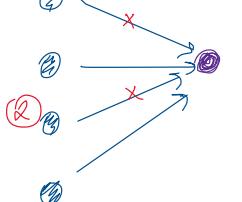
$\alpha_3$  / = keep prob (Inverted dropout technique)

⑪ (don't dropout at test time.)

# Understanding Dropout

(#) Can't rely on one feature, so have to spread weights  
≈ Share weights

~~(A) 10~~



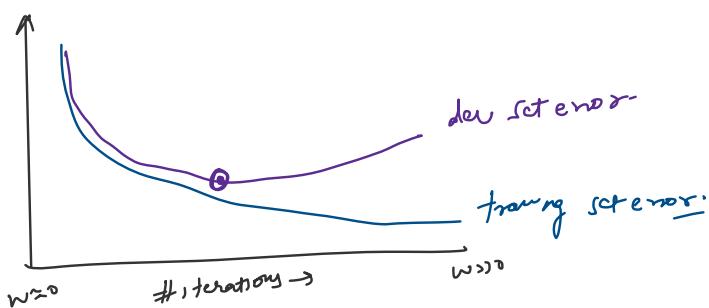
④ Can vary, keep prop according to layers,

$$\# \text{ units} \propto \frac{1}{\text{Keep Prop Value}}$$

## Computer Vision

## Other Regularization Methods

## EARLY STOPPING



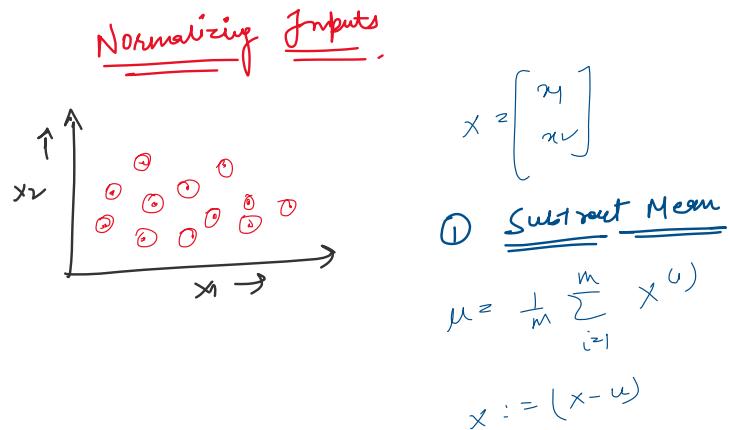
# optimize cost fx or (Gradient descent, LBGK, ...)

# do not overfit  
(Regularization, getting more data  $\Rightarrow$ )

Orthogonalization → Thinking of one task at a time.

- Orthogonalization → Thinking of one task at a time  
 → ① "no orthogonalization is early stopping"
- $L_2$  Regularization

## Setting up your Neural N/W



② Normalize Variances

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m [x^{(i)} - \mu]^2 \quad \text{(Squaring.)}$$

$$x = \frac{x}{\sigma}$$

$$x = \frac{(x - \mu)}{\sigma}$$

③ use same  $\mu$  and  $\sigma^2$  to normalize Test set.

Why NORMALIZE?

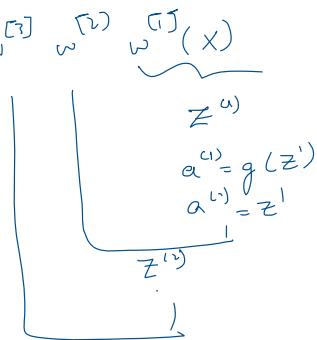
$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})$$

$$\begin{aligned} w_1 &= 1 && \dots & 1000 \\ w_2 &= 1 && \dots & 10 \end{aligned}$$

Variance | Exploding gradients



$$\hat{y} = w^{[1]} + w^{[2]} * g(z) = \hat{z}$$



$$W^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \quad \text{with } \lambda^{(l)} = 10^5$$

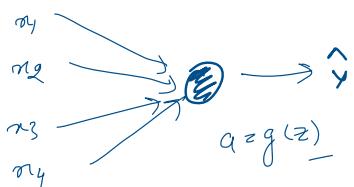
$$\hat{y} = W^{[l]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{-1} \times \hat{y}$$

$\hat{y} = (1.5)^L$

- $\rightarrow 10^5$  times number of layers
- $\rightarrow$  if you have deep neural net, value of  $\hat{y}$  will explode.

### Weight Initialization For Neural N/W.

#### # Single Neuron



$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$b = 0$$

$$n \propto \frac{1}{w_i}$$

$$\text{Var}(w_i) = \frac{1}{n} \quad [n \in \text{No. of input features in neuron}]$$

$$W^{[l]} = np.random.randint(\text{shape}) * np.sqrt\left(\frac{1}{n^{[l-1]}}\right)$$

$2/n$  for ReLU ( $f(x)$ )

$$g^{(l)}(z) \in \text{ReLU}$$

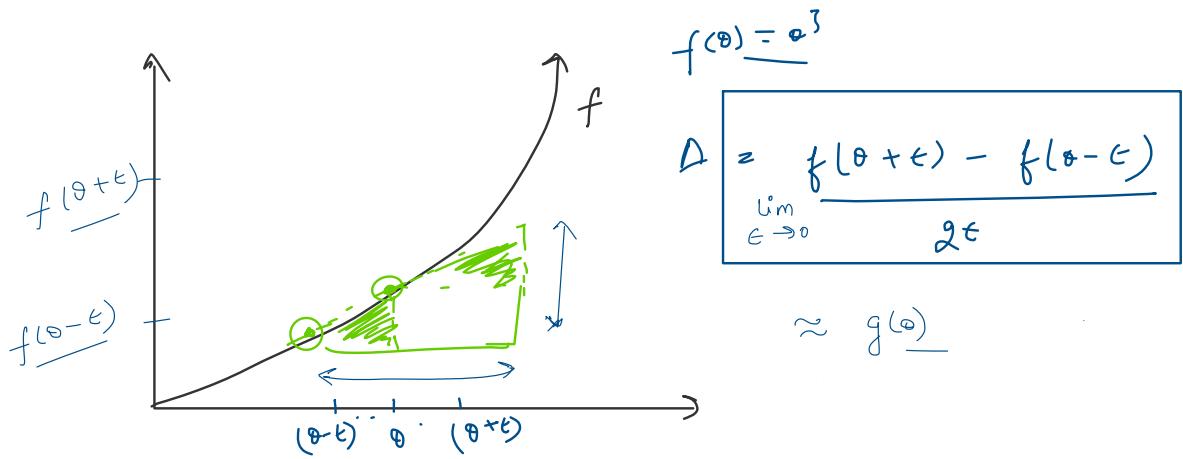
#### Other Variants

$$\tanh(z) = \sqrt{1/n^{[l-1]}} \quad \left. \begin{array}{l} \text{xavier} \\ \text{initialization} \end{array} \right\}$$

### Numerical Approximation Of Gradients

~ Backprop (Gradient checking) -

$$\uparrow \quad \uparrow_n \quad f(\theta) = 0^3$$



### Gradient Checking.

gradient checking for Neural N/W  $\rightarrow$

Take  $w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}$  & reshape to a big vector,  $\underline{\theta}$ ,  
Concatenate

Take  $d\theta^{(1)}, d\theta^{(2)}, d\theta^{(3)}, \dots, d\theta^{(L)}$  & reshape into a big vector,  $d\underline{\theta}$   
concatenate  $\rightarrow$

grad check

for each  $i \rightarrow$

$$d\theta_{approx}^{(i)} := \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon)}{2\epsilon}$$

$$\approx d\theta^{(i)}$$

$$= \frac{\partial J}{\partial \theta_i}$$

$$\partial \theta_{approx} \approx \partial \theta$$

$$\text{check} = \frac{\|\partial \theta_{approx} - \partial \theta\|_2}{\|\partial \theta_{approx}\| + \|\partial \theta\|_2}$$

$$\epsilon \approx 10^{-7}$$

### gradient Checking Implementation

- ④ don't use grad check in training only to debug,

\* if algorithm fails in grad check; look at the components & try to find the bug.

$$\partial b^{[l]}, \partial w^{[l]}$$

\* Remember regularization

$$J(\theta) = \frac{1}{m} \sum L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_2^2$$

$$\partial \theta = \text{grad of } J \text{ wrt. } \theta$$

\* grad check doesn't work with dropout

\* run at random initialization; perhaps after some training

Mini batch gradient descent

$$\mathbf{X} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \dots, \mathbf{x}^{(m)}] \\ (n \times m)$$

$$\mathbf{Y} = [y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(m)}] \\ (1, m)$$

if training examples are very large, say 10 million.

① will have to traverse (10 million) training set before you take one little step of gradient descent.

Mini-Batch

Now, One mini batch consists of

$$\text{if } X=50 \text{ (say)} \text{ & } n=5$$

$$X^{\{1\}} = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & x^{(4)} & \dots & x^{(5)} \end{bmatrix}$$

$$X^{\{1\}} = \dots$$

$$\frac{\text{training-examples}}{\text{n. of mini-batches}} = n(\text{mb}).$$

→ and same for Y.

# Mini-Batch t: $X^{t+3}, Y^{t+3}$	① $X_{\text{dim}} = n_x, n(\text{mb})$
	② $Y_{\text{dim}} = 1, n(\text{mb})$

Mini-Batch gradient descent

for  $t = 1 \dots n(\text{mb})$ :

// perform a step of gradient descent using  $X^{\{t\}}, Y^{\{t\}}$

① Perform forward prop.  $\Rightarrow$

$$Z^{[0]} = W^{[0]} X^{[0]} + b^{[0]}$$

$$A^{[0]} = g^{[0]}(Z^{[0]})$$

$$A^{[L]} = g^{[L]}(Z^{[L]})$$

② Compute loss

$$J = \frac{1}{m} \sum J(\rightarrow) + \text{regularization term}$$

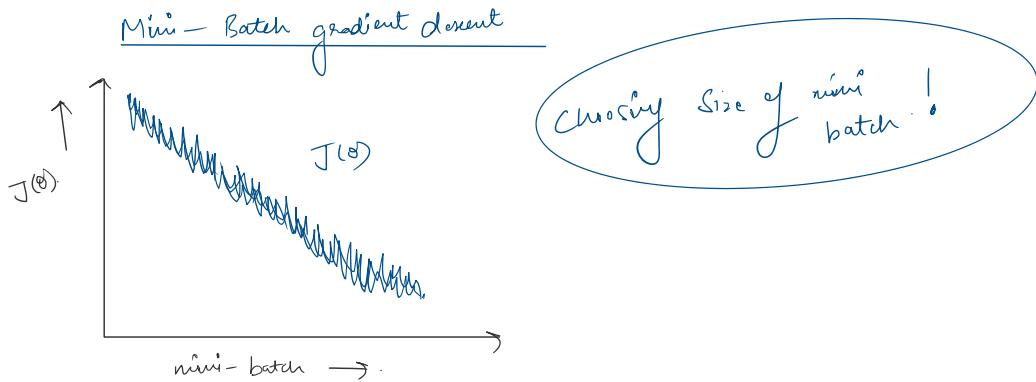
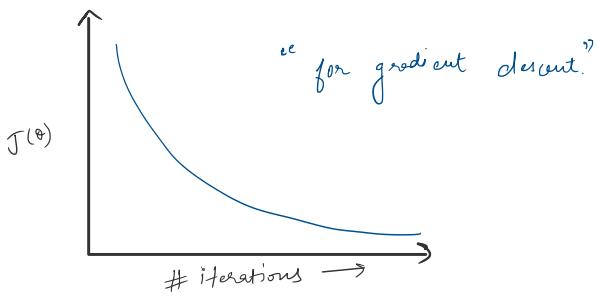
$$W^{[L]} := W^{[L]} - \alpha \partial W^{[L]} \quad \& \text{ similarly,}$$

$$b^{[L]} := b^{[L]} - \alpha \partial b^{[L]}$$

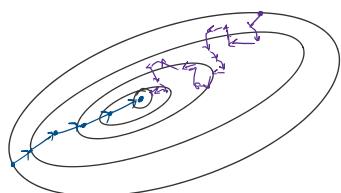
"Also called 1 epoch / 1 pass"

L2. UNDERSTANDING MINI BATCH GRADIENT DESCENT.

## L2. UNDERSTANDING MINI BATCH GRADIENT DESCENT.



- ④ if mini batch size =  $m$ , then : Batch gradient descent.
- ⑤ if (mini batch size = 1) : Stochastic gradient descent.  
 $(x^{(1)}, y^{(1)}, x^{(2)}, y^{(2)}, \dots)$



[In practice mini batch size is between 1 and  $m$ .]

### Batch gradient descent

mini-batch size :  $m$

↓  
Too long per iteration

### Stochastic gradient descent

batch size : 1

↓  
loop speeding up from vectorization.

⑥ mini - batch gradient descent

### Choosing size of mini-batch

if Small training set : Use batch gradient descent.

if Small training set : Use batch gradient descent.  
( $m \leq 2000$ )

Typical mini-batch sizes :  $2^k$  (multiple of 2.)

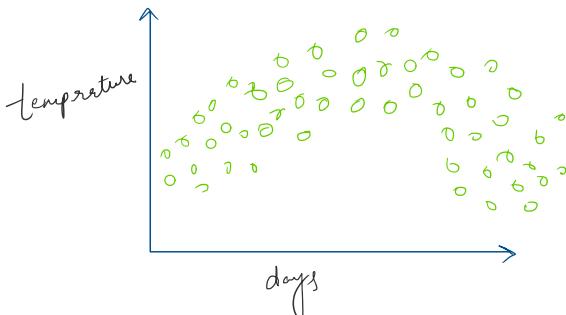
Make sure mini-batch fits in CPU / GPU memory.

$X^{(t)}$ ,  $y^{(t)}$  → Hyperparameter.

## L3 EXPONENTIALLY WEIGHTED AVERAGE

Temperature in New-delhi →

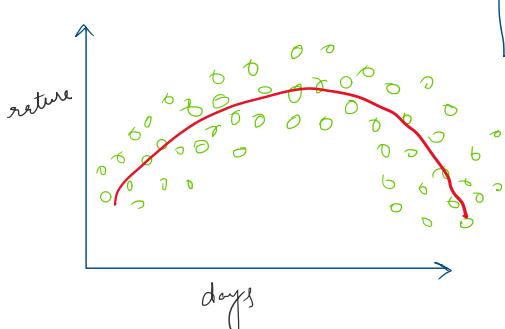
$$\begin{aligned} \theta_1 &= 32^\circ C \\ \theta_2 &= 31^\circ C \\ \theta_3 &= 29^\circ C \\ &\vdots \\ \theta_{190} &= 28^\circ C \end{aligned}$$



$$\begin{aligned} V_0 &= 0 \\ V_1 &= 0.9V_0 + 0.1\theta_1 \\ V_2 &= 0.9V_1 + 0.1\theta_2 \\ V_3 &= 0.9V_2 + 0.1\theta_3. \end{aligned}$$

$$V_t = 0.9V_{(t-1)} + 0.1\theta_t$$

→ Exponentially weighted average.  
→ Moving Average.



$$V_t = \beta V_{(t-1)} + (1-\beta)C$$

Exponentially weighted moving average.

generalize formula.

$$\beta = 0.9$$

$$V_t \stackrel{(avg)}{=} \frac{1}{(1-\beta)} \text{ days } \text{ temperature.}$$

L4

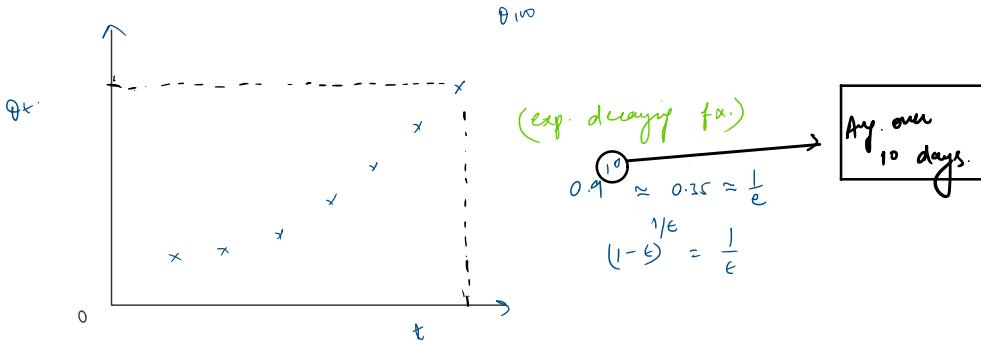
## UNDERSTANDING EXPONENTIALLY WEIGHTED Avg.

$$V_t = \beta V_{(t-1)} + (1-\beta) \theta_t$$

$$V_{100} = 0.1 \theta_{100} + 0.9 V_{99} \quad \beta = 0.9$$

$$= 0.1 \theta_{100} + 0.1 * 0.9 (\theta_{99}) + 0.1 (0.9)^2 \theta_{98} + \dots$$

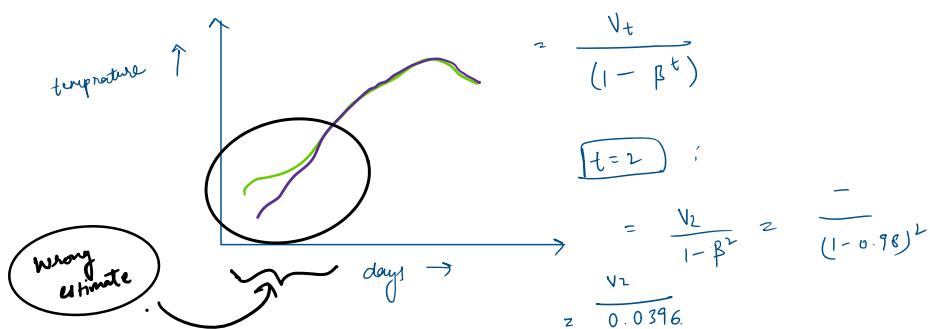
- exponentially decaying function.



$\cdot V_0 = 0$ $\cdot V_t = \beta V_{(t-1)} + (1-\beta) \theta_t$	$V = 0$ $V_0 = \beta V + (1-\beta) \theta_1$ $V_1 = \beta V + (1-\beta) \theta_2$
---	---

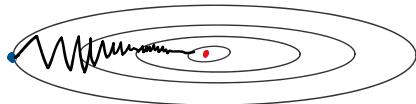
## L5 BIAS CORRECTION IN EXP. WEIGHTED AVG.

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t \quad ; \quad \beta = 0.98$$



## L6 GRADIENT DESCENT WITH MOMENTUM

↑ - slower learning.  
 ↘ - faster learning.



# Momentum :-

On every iteration  $t$ , calculate  $\nabla w$

$\nabla w, \nabla b$  on current batch

friction

$$\nabla w = \beta \nabla w + (1-\beta) \nabla w \quad ; \quad \text{Acceleration}$$

friction

$\delta w, \delta b$  on current batch

Velocity

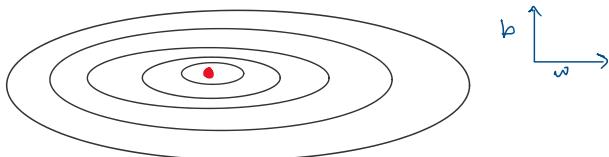
$$\begin{aligned} V_{\delta w} &= \beta V_{\delta w} + (1-\beta) \delta w \\ V_{\delta b} &= \beta V_{\delta b} + (1-\beta) \delta b \end{aligned} \quad \left. \begin{array}{l} \text{Acceleration} \\ \vdots \end{array} \right.$$

$$w := w - \alpha V_{\delta w}$$

$$b := b - \alpha V_{\delta b}$$

Hyperparameters :-  $\alpha, \beta$        $\boxed{\beta = 0.9}$ .

## L7. RMS Prop (Root mean square prop)



On every iteration, t :-

Compute  $\delta w, \delta b$  for mini-batch.

$$S_{\delta w} = \beta S_{\delta w} + (1-\beta) \delta w^2$$

element wise square  $\Rightarrow$   $p$  is hyperparam for momentum

$$S_{\delta b} = \beta S_{\delta b} + (1-\beta) \delta b^2$$

$$\boxed{S_{\delta w} \approx}$$

to make sure it's tve, adding small value  $\epsilon$ .  
 $\epsilon \approx 10^{-8}$

$$w := w - \frac{\alpha \delta w}{\sqrt{S_{\delta w}} + \epsilon}$$

①  $w \rightarrow$  : fast.  
 $b \uparrow$  : slow down.

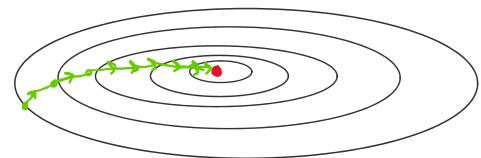
$$b := b - \frac{\alpha \delta b}{\sqrt{S_{\delta b}} + \epsilon}$$

$S_{\delta w}$  = relatively small.

$S_{\delta b}$  = relatively large.

\* RMS prop.

→ now, can also use large learning rate,  $\alpha$ .



## L8. ADAM'S OPTIMIZATION ALGORITHM

$$\{ \text{grad (Momentum)} + \text{Rmsprop.} \}$$

$$V_{\delta w} = 0, \quad S_{\delta w} = 0 ; \quad V_{\delta b} = 0, \quad S_{\delta b} = 0$$

On every iteration t :-

\* Compute  $\delta w, \delta b$  using current mini-batch.

$$\begin{aligned} V_{\delta w} &= \beta_1 V_{\delta w} + (1-\beta_1) \delta w \\ V_{\delta b} &= \beta_1 V_{\delta b} + (1-\beta_1) \delta b \end{aligned} \quad \left. \begin{array}{l} \text{• Momentum (Hyp: } \beta_1 \text{)} \\ \vdots \end{array} \right.$$

$$\begin{aligned} S_{\delta w} &= \beta_2 S_{\delta w} + (1-\beta_2) \delta w^2 \\ S_{\delta b} &= \beta_2 S_{\delta b} + (1-\beta_2) \delta b^2 \end{aligned} \quad \left. \begin{array}{l} \text{• rms prop.} \\ (\text{Hyp: } \beta_2) \end{array} \right.$$

### Bias correction

$$* \quad \begin{aligned} V_{dw}^{\text{corrected}} &= V_{dw} / (1 - \beta_1^t) \\ V_{db}^{\text{corrected}} &= V_{db} / (1 - \beta_1^t) \end{aligned} \quad \left. \right\} \cdot \text{Momentum}$$

$$\left. \begin{aligned} S_{dw}^{\text{corrected}} &= S_{dw} / (1 - \beta_2^t) \\ S_{db}^{\text{corrected}} &= S_{db} / (1 - \beta_2^t) \end{aligned} \right\} \cdot \text{Rms Prop.}$$

$$* \quad w_i = w_i - \alpha \left[ \frac{V_{dw}^{\text{corrected}}}{(\sqrt{S_{dw}^{\text{corrected}}} + \epsilon)} \right]$$

$$* \quad b_i = b_i - \alpha \left[ \frac{V_{db}^{\text{corrected}}}{(\sqrt{S_{db}^{\text{corrected}}} + \epsilon)} \right]$$

### hyperparameters choice :-

$\alpha$  = Needs to be tuned.

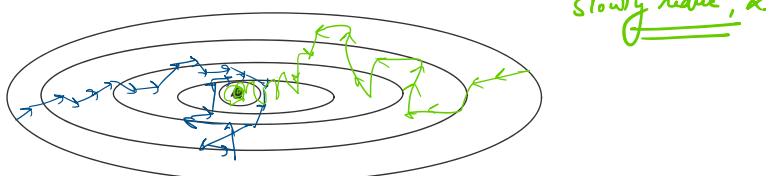
$\beta_1 = 0.9$  (dw)

$\beta_2 = 0.999$  ( $\delta w^2, \delta b^2$ )

$\epsilon = 10^{-8}$

ADAM : "Adaptive moment Estimation."

### L9. LEARNING RATE DECAY:



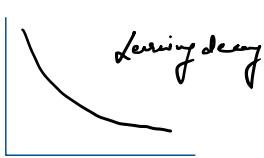
\* 1 epoch = 1 pass over the data.

$$\alpha = \left( \frac{1}{1 + \text{decay rate} * \text{epochnum}} \right) \alpha_0 \quad \xrightarrow{\text{some initial learning rate.}}$$

Epoch	$\alpha$
1	0.1
2	0.067
3	0.05
4	0.04

$$\alpha_0 = 0.2$$

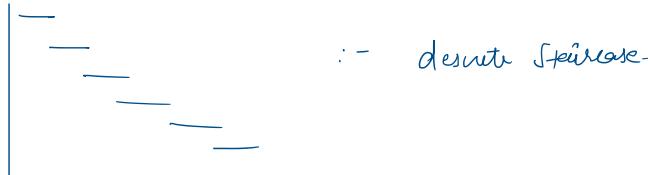
$$d\tau = 1$$



$$\alpha = \left[ (N < 1) \text{ epochnum.} \dots \alpha_0 \right] \vdash \text{exponential decay.}$$

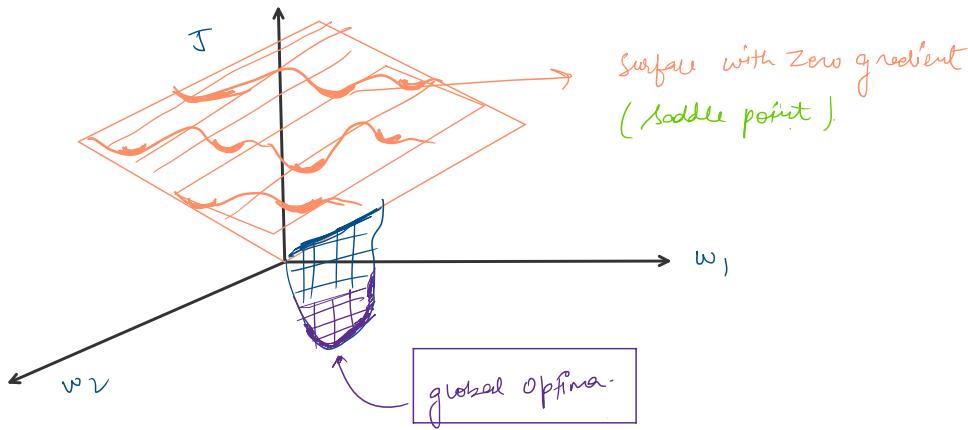
people also use:

$$\alpha = \frac{k}{\sqrt{\text{epoch}}} * \alpha_0$$



\* Manual decay ..

## L10: THE PROBLEM OF LOCAL OPTIMA



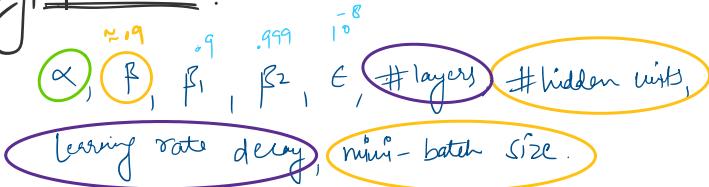
problem with plateaus.

- \* [ plateau is a region where the derivative is 0 for a very long time. ]
  - Unlikely to stuck in a bad local optima.
  - Plateaus can make learning slower.

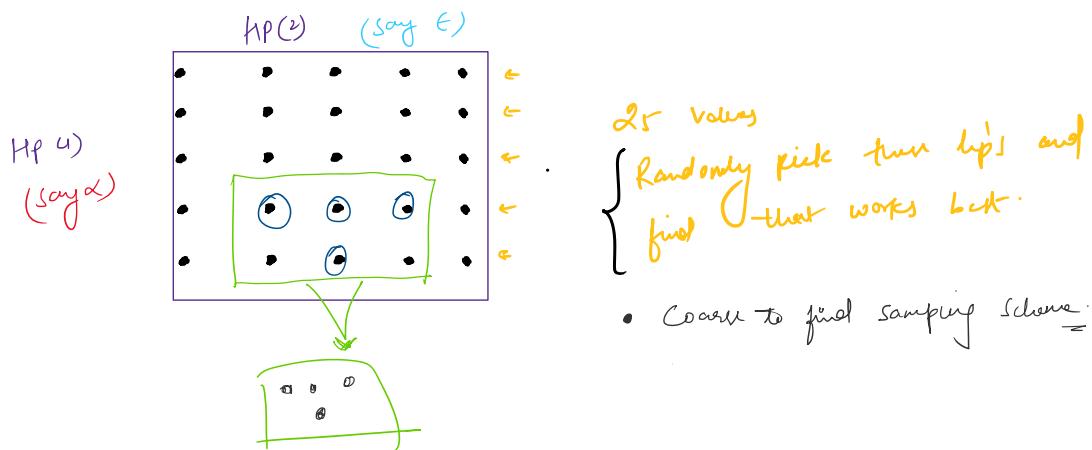
## Hyperparameter Tuning

### L1. Tuning process.

#### Hyperparameters



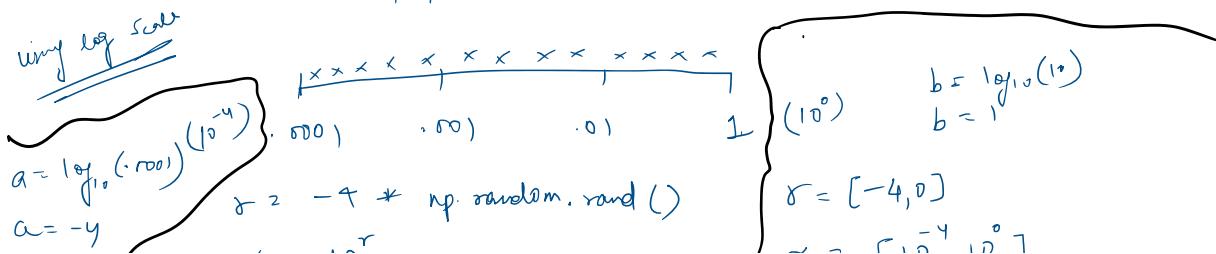
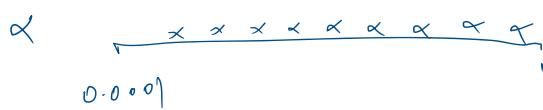
- \* - high priority.
- \* - medium priority.
- \* - low priority.
- \* - almost negligible (default values are used.)



\* Use random Sampling and adequate search.

{optional} \* Coarse to fine search process.

### L2. Using An Appropriate Scale To Pick Hps.



$a = \log(1 - \text{rand})$   
 $a = -y$   
 $\alpha = 10^r$   
 Sample  $r$  uniformly on  $\in [a, b]$

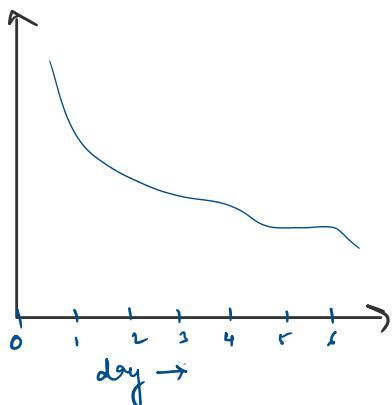
$\gamma = [-4, 0]$   
 $\alpha = [10^{-4}, 10^0]$

### # Hints for exponentially weighted avg.

$\beta = [0.9, \dots, 0.999]$   
 Linear Scale —  $\times$   
 $1 - \beta = [0.1, \dots, 0.001]$   
 1 → .1 → .01 → .001  
 ✓  
 very sensitive to small changes in  $\beta$ .

### L3: HP tuning (Panda vs Caviar)

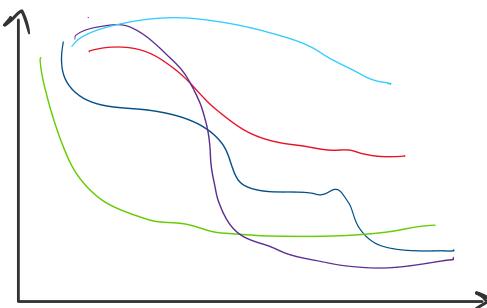
#### Babysitting One Model



Keep checking one day at a time.

#### Panda Approach

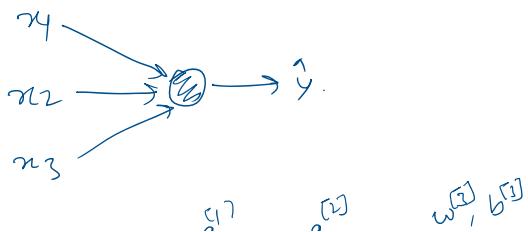
#### Training many models in Parallel



#### Caviar Approach

### Batch Normalization

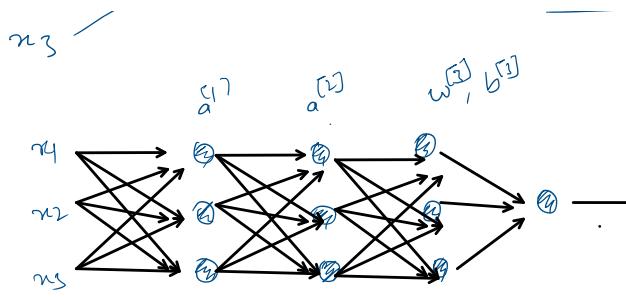
#### Normalizing Activations in a NW



$$x = x - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_i x_i^2$$

$$x / \sigma^2$$



Can we normalize  $a^{[2]}$ ? to train  $w^{[3]}, b^{[3]}$  faster?

\* Normalize  $z^{[2]}$  or  $z^{[3]}$  instead of  $a^{[2]}$

$$\mu = \frac{1}{m} \sum z^{[i]}$$

$$\sigma^2 = \frac{1}{m} \sum (z^{[i]} - \mu)^2$$

$$z_{\text{normalize}}^{[i]} = \frac{z^{[i]} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$z$  has normal distribution and variance 1, but we want different dist.

$$\tilde{z}^{[i]} = \gamma z_{\text{norm}}^{[i]} + \beta$$

Learnable parameter.

Standardize mean  
and Variance.

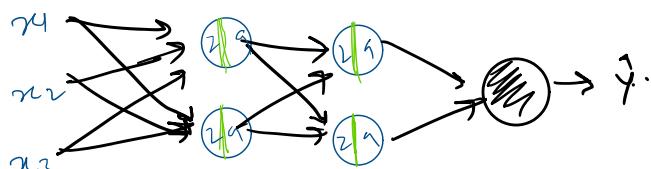
$$\gamma = \sqrt{\sigma^2 + \epsilon} \quad \text{and} \quad \beta = \mu$$

then:

$$\tilde{z}^{[i]} = \gamma z^{[i]} + \beta$$

Use  $\tilde{z}^{[i]}$  instead of  $z^{[i]}$

## L2. Fitting Batch Norm to a N/W $\Rightarrow$



$$x \xrightarrow{w^{[0]}, b^{[0]}} z^{[1]} \xrightarrow{\gamma^{[1]}, \beta^{[1]}, \mu^{[1]}, \sigma^2} \tilde{z}^{[1]} \xrightarrow{a^{[1]} = g(\tilde{z}^{[1]})} z^{[2]} \dots$$

Parameters  $\Rightarrow w^{[0]}, b^{[0]}, \beta^{[1]}, \gamma^{[1]}, \mu^{[1]}, \sigma^2$

$$\beta^{[1]} := \beta^{[1]} - \alpha \cdot d\beta^{[1]}$$

### Working with mini-batches.

$$x^{[1]} \xrightarrow{w^{[0]}, b^{[0]}} \tilde{z}^{[1]} \xrightarrow{\gamma^{[1]}, \beta^{[1]}, \mu^{[1]}, \sigma^2} g^{[1]}(\tilde{z}^{[1]}) \xrightarrow{a^{[1]}} \dots$$

$$\hat{x} = \overbrace{w^{[0]} b^{[0]}}^{\text{forward pass}} \xrightarrow{(f \cdot n)} \underbrace{0}_{\text{batch norm}} \xrightarrow{\gamma^{[0]} \beta^{[0]}} \hat{x}$$

$$\text{Parameters} = w^{[0]}, b^{[0]}, \beta^{[0]}, \gamma^{[0]}, \Rightarrow \hat{x} = w^{[0]} a + b^{[0]}$$

\* Batch Norm averages all  $\hat{x}^{[0]}$ , so adding/subtracting any const. doesn't affect the equation.

$$\hat{x}^{[0]} = w^{[0]} a^{[0]}$$

$$\hat{x}_{\text{norm.}}^{[0]} = \frac{\hat{x}^{[0]} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{x}^{[0]} = \gamma^{[0]} \hat{x}^{[0]} + \beta^{[0]}$$

$$\text{dim}_a := (n^{[0]}, 1)$$

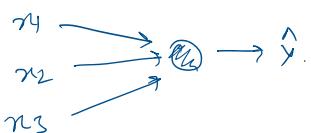
Implementing gradient descent.

for  $t=1$  to ... mini-Batch<sub>n</sub> :-

- ① Compute forward prop. on  $x^{[t]}$
- ② use batch norm to replace  $\hat{x}^{[0]} \Rightarrow \tilde{x}^{[0]}$
- ③ use backprop to calculate grads.
- ④ update parameters,  $\rightarrow w^{[0]} := w^{[0]} - \alpha \delta w^{[0]}$   
 $\rightarrow \beta^{[0]} := \beta^{[0]} - \alpha \delta \beta^{[0]}$   
 $\rightarrow \gamma^{[0]} := \gamma^{[0]} - \alpha \delta \gamma^{[0]}$

### L3. Why does batch norm works?

\* Leaving on shifting input dist.

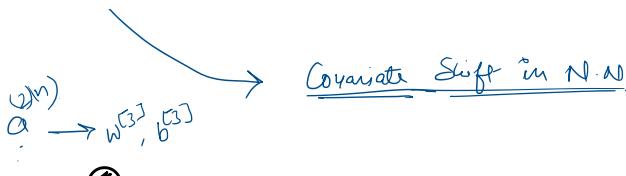


Trained on monochrome image

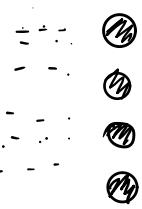
Trained on coloured image

"Covariate Shift"  $\rightarrow$

\* If you learn some  $x \rightarrow y$  mapping, if distribution of  $x$  changes then you need to retrain your learning algorithm.



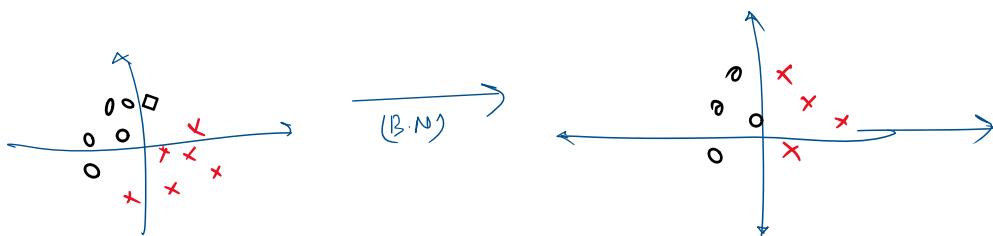
Covariate Shift in N.W.



~~Problem~~

parameters  $w, b$  before this layer change,  $a^{(2)}$  will also change.  
Suffering from a problem of "covariate shift".

- # Batch norm : No matter the distribution changes, mean and variance remains same.



## # Batch Norm as Regularization.

- ① Each mini-batch is scaled by mean/variance computed on just one batch.
- ② This adds some noise  $\tilde{z}^{(i)}$ , scaling process of  $z^{(i)} \rightarrow \tilde{z}^{(i)}$  becomes more noisier, similar to dropout.
- ③ This has slight regularization effect.

mini-batch size  $\propto \frac{1}{\text{noise}} \propto \text{regularization}$ .



"don't use batch norm as regularization  
only use for normalization."



## L4: Batch Norm at -last time

We need different  $\mu$  and  $\sigma^2$ .  
|

$$\mu = \frac{1}{m} \sum z^{(i)}$$

$m$  :- ex: in mini-batches

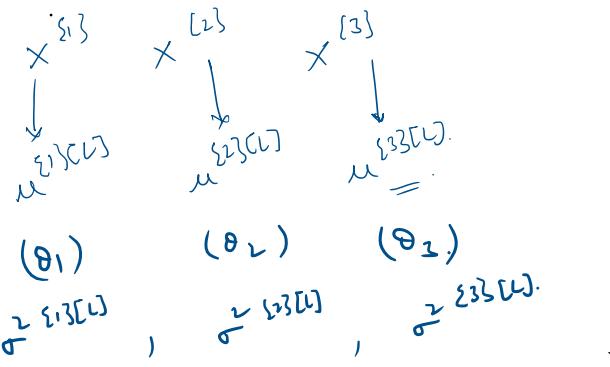
$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

$$\rightarrow z^{(i)} - \frac{1}{m} \sum z^{(i)} - \mu$$

$$Z_{\text{norm}}^{(i)} = \frac{Z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{Z}^{(i)} = \gamma Z_{\text{norm}} + \beta$$

estimate using exponentially weighted avg. (across mini-batch)



new  
Znorm

$$\tilde{Z}_{\text{norm}} = \frac{\tilde{Z} - \mu}{\sqrt{\sigma^2 + \epsilon}}, \quad \tilde{Z} = \left[ \gamma Z_{\text{norm}} + \beta \right].$$

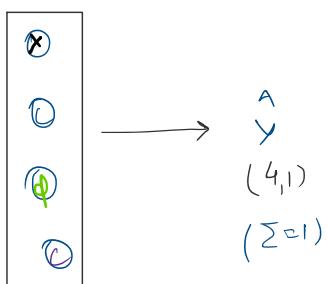
### Multi Class Classification

#### Softmax Regression.

Now recognize cats, dogs and chicks.

Class 1    Class 2    Class 3    -1

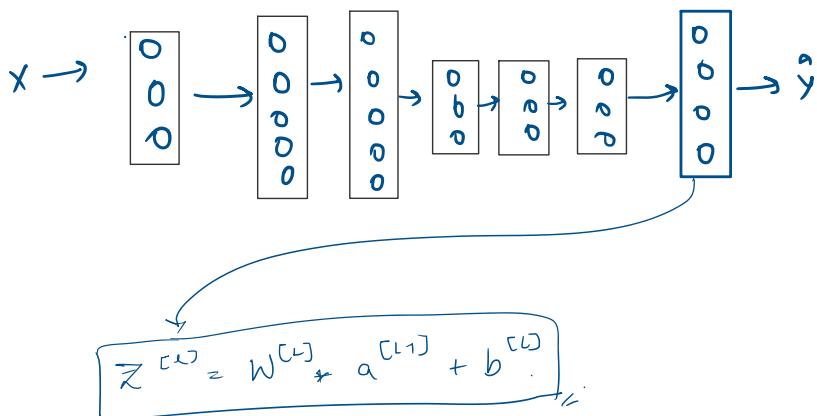
# classes : 4



- final layer -

Softmax Layer

## Softmax Layer



## Softmax Activation fn.

$$t = e^{(\tilde{z}^{[L]})}$$

$$(4,1) \quad a^{[L]} = \frac{e^{(\tilde{z}^{[L]})}}{\sum t^i}, \quad a_i^{[L]} = \frac{t^i}{\sum_{i=1}^4 t^i}$$

$$\tilde{z}^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$\sum_{i=1}^4 t^i = 176.3$$

$$a^{[L]} = \frac{t}{176.3} = \frac{1}{176.3} * \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$a^{[L]} = g^{(L)}(\tilde{z}^{[L]}).$$

$\uparrow_{4 \times 1}$

L2:

Training a softmax classifier:

$$\tilde{z}^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}, \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$a^{(v)} = g^{(v)}(z^{(v)}) = \begin{pmatrix} e^5 / \sum e \\ e^2 / \sum e \\ e^1 / \sum e \\ e^3 / \sum e \end{pmatrix} = \begin{pmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{pmatrix}$$

"hard max"

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{Class with highest probability.}$$

- \* Softmax regression generalizes Logistic Regression to  $c$  classes.  
if  $c=2$ , Softmax Regression = Logistic Regression.

### Loss fx.

$$y^{(x_i)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \xrightarrow{\text{Cat.}} \hat{y}^{(x_i)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$L(\hat{y}, y) = - \sum_{j=1}^c y_j (\log \hat{y}_j)$$

$$(y_1 = y_3 = y_4 = 0) : 3 \text{ values in } L(\hat{y}, y) = 0$$

$$L(\hat{y}, y) = -y_2 \log \hat{y}_2$$

(minimize)  $\hat{y}_2 = 1 \Rightarrow -\log \hat{y}_2$

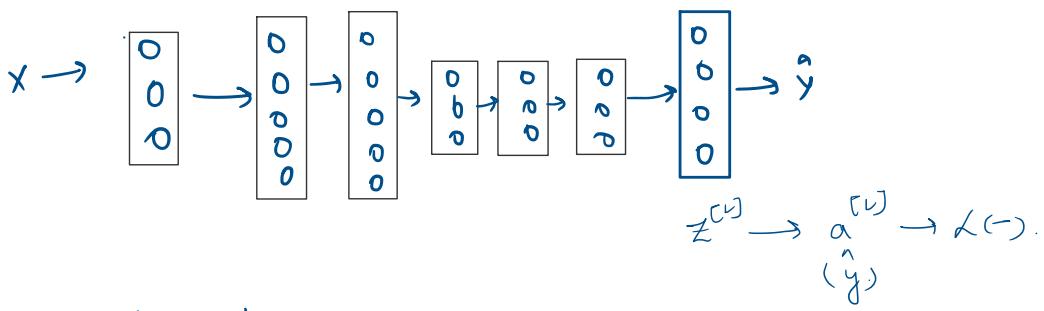
make  $\hat{y}$  as big

$$J(\omega^{(v)}, b^{(v)}) = 1/m \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$y = [y^{(1)} \ y^{(2)} \ y^{(3)} \dots \ y^{(m)}]_{(4 \times m)}$$

$$\hat{y} = [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \hat{y}^{(3)} \dots \ \hat{y}^{(m)}]_{(4 \times m)}$$

### Gradient Descent with Softmax.



forward prop : ✓

backward prop :-

$$\delta z^{[L]} = \hat{y} - y \quad (\in \mathbb{R}^n) \quad (\leftarrow)$$

$\frac{\partial J}{\partial z^{[L]}}$

### Deep Learning Frameworks

- Caffe | Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mx net
- PaddlePaddle
- PyTorch
- TensorFlow
- Theano

## Week 1 - ML strategy (1)

13 June 2020 11:14

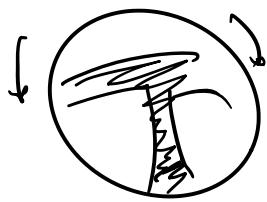
### Introduction to ML strategy

#### L2. Orthogonalization

\* Tuning one property at a time.

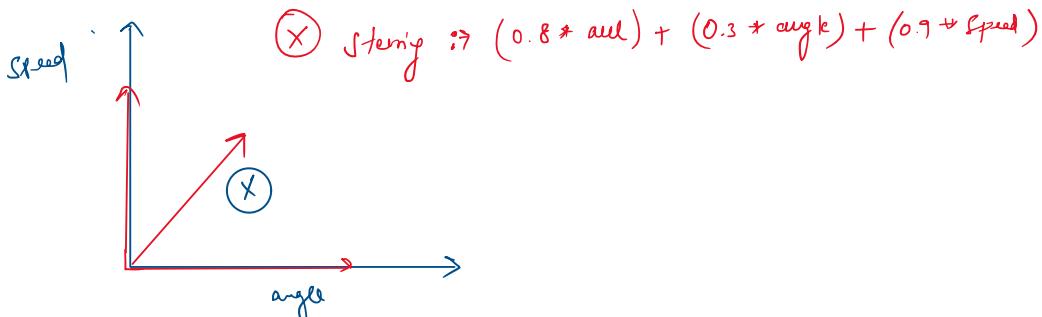
ex:

driving a car



- direction  $\Rightarrow$  steering
- Acceleration, • Braking

All have separate functions to tune



#### Chain of Assumptions in ML

\* fit training set well on cost function.



↓  
dev set

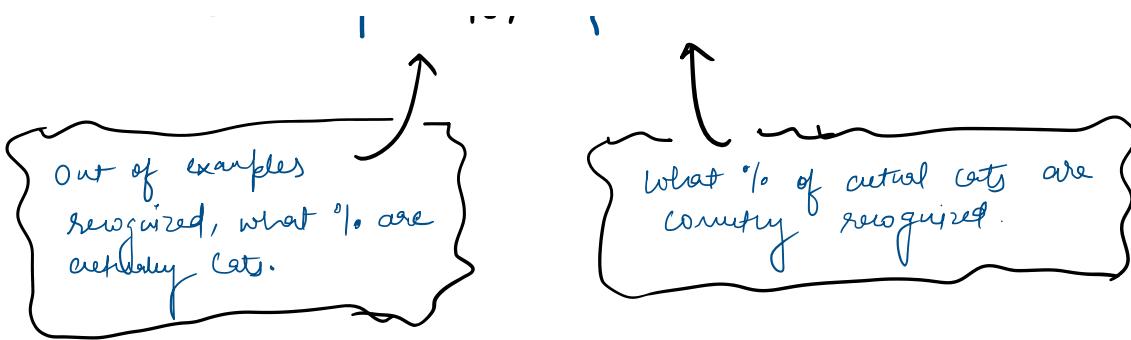


### Setting up your goal

#### L1. Single Number evaluation metric

Classifier	Precision	Recall
A	95%	90%
B	98%	85%

↑                      ↑



## What to choose?

$$\underline{\underline{F_1\text{-score}}} \quad \text{"Avg. of Precision and recall"} \\ \left( \frac{2}{1/P + 1/R} \right) \quad \therefore \text{Harmonic Mean}$$

L2. Satisfying and optimizing metrics.

<u>Classifier</u>	<u>Accuracy</u>	<u>Running Time</u>
A	90%	80 ms.
B	92%	95 ms.
C	95%	1500 ms.

- maximize A
- minimize Running Time.

N metrics :- 1 Optimizing metric.  
 $(N-1)$  satisfying metrics.

L3. Train | dev | Test distributions.

## Cat Classifier.

~~Regions~~  $\Rightarrow$

dev			
• India	• Mexico	• Canada	
• Spain	• Japan	• China	

~~test~~



# Dev/Test and test-set should come from same distribution.

#### 14- Size of Dev/test Sets

Train :- 70%

test :- 30%

or

Train :- 60% dev/cv :- 20% test :- 20%

for 1000, nos examples -



#### Size of Test Set

- \* Set your test size big enough to give high confidence in the overall performance of system.

#### 15. When to change dev/test sets and metrics

Metric :- Classification error.

- (+) • Algorithm A :- 3% error → photographic images
- (✓) • Algorithm B :- 5% error. → ! p-hyp. }  
} Misclassification.

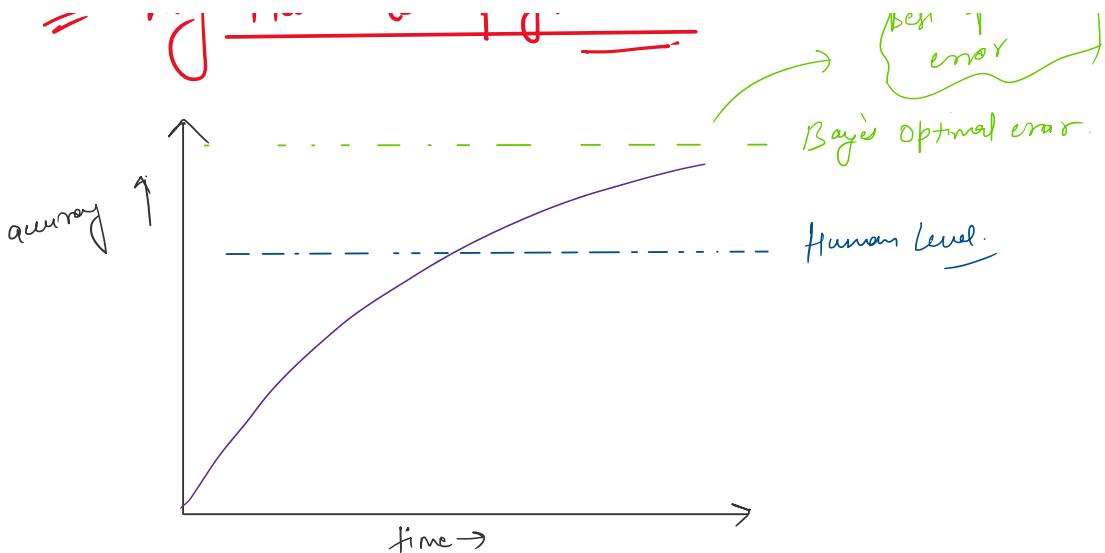
$$\text{Error} \Rightarrow \frac{1}{\sum w^i} \sum_{i=1}^{m_{\text{data}}} \{ y_{\text{pred}}^{(i)} \neq y^{(i)} \}$$

$$w^i = \begin{cases} 1, & \text{if correct output} \\ 100, & \text{if misclassified.} \end{cases}$$

Comparing to human level performance

#### 16. Why Human level performance

best optimal error



### Avoidable Bias

Human Error  $\rightarrow 1\%$

Training error  $\rightarrow 8\%$

dev Error  $\rightarrow 10\%$

7.5% (say)

8%

10%

Avoidable bias

\* focus on bias

\* focus on variance

# Human level error is a proxy for Bayes error.

almost same content as in  
previous course  
or  
Machine Learning by  
Andrew Ng.

## Error Analysis

II. Carrying Out error Analysis.

Cat classifier - 5% error / 95% accuracy.

- Get 100 ~ mislabelled dev set examples.
- Count up how many are dogs.
- Count errors.

Bias and Variance with mismatched data distribution

Assume Human got  $\approx 1\%$  error.

Training error  $\Rightarrow 1\%$ .

Dev Error  $\approx 10\%$ .

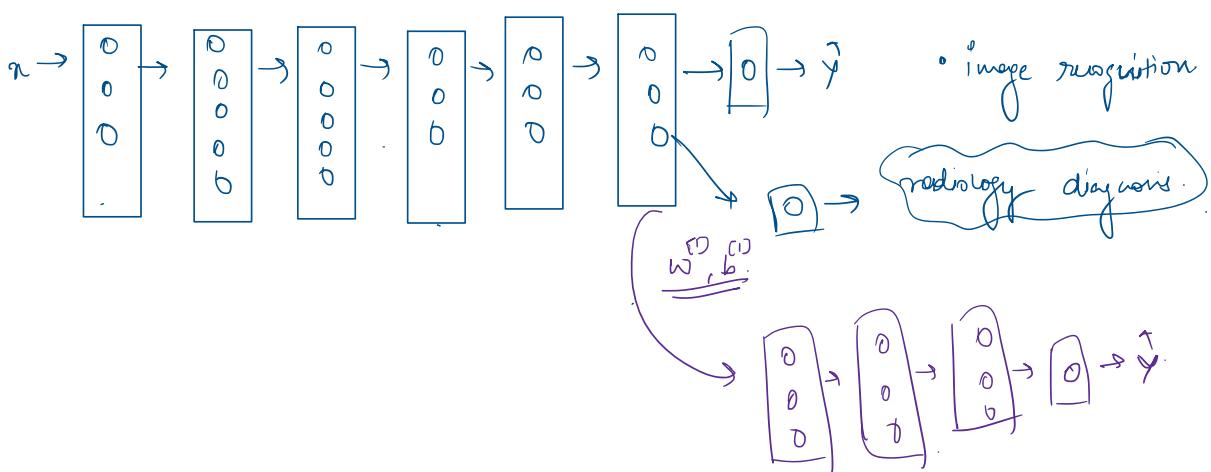


Variance Error.

# Explain almost same as train/test | dev split in 2nd Course.

Date mismatch problem:

- Carry out manual analysis and try to understand the difference b/w training and dev/test sets.
- Make train set more similar.

Learn from Multiple Task.Transfer Learning  $\Rightarrow$ 

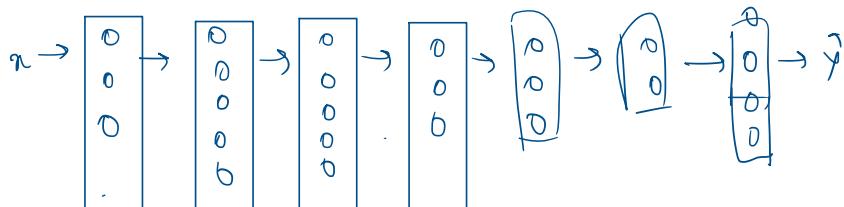
When transfer learning makes sense  $\Rightarrow$

When transfer learning makes sense  $\Rightarrow$

Transfer(A)  $\rightarrow$  (B) { They have same input  $x$ , A has more data than B }

### Multi-Task Learning.

\* Start off simultaneously



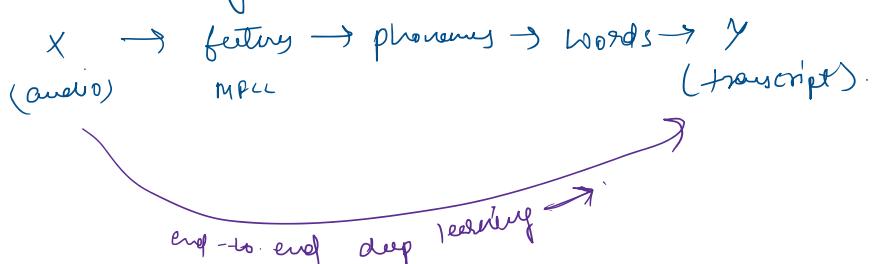
$$\text{Loss } \hat{y}^{(i)} = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^q L(\hat{y}^{(i)}, y_j)$$

(Logistic Loss.)

\* Trains a set of tasks that could benefit from sharing shared low level features.

### End-to-end deep learning.

Speech Recognition (ex.)



### Whether to use End-to-end dep. learning

#### Pros

- Let the data speak.
- Few hand design components needed.

#### Cons

- Needs large amount of data.  
 $X \rightarrow Y$
- Excludes potentially useful hand design components.

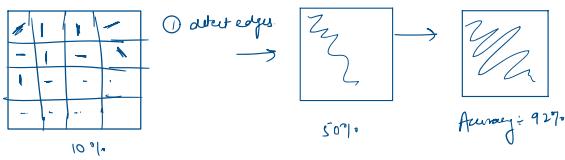
## Convolutions Neural Networks

Computer Vision

- Object detection.
- Image classification.
- Neural style transfer.

  $64 \times 64 := (64 \times 64) \times 3$   
 $\{rgb\}$

  $1024 \times 1024 := (1024 \times 1024) \times 3$   
 $\approx 3 \text{ million.}$   
(Very high)

Edge DetectionComputer Vision problems

- ① detect vertical lines: (|)  
② detect horizontal lines: (—)

Vertical edge detection

\* (Shift up 1 column after iteration k row after iteration j row)

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

Image:  $(6 \times 6)_{\times 1}$

Convolution:

1	0	1
1	0	-1
1	0	-1

$= 4 \times 4 \text{ mat.}$

filter / kernel.

$(3 \times 1) + (0 \times 0) + (1 \times -1) + (1 \times 1) + (5 \times 0) + (8 \times -1) + (2 \times 1) + (7 \times 0) + (2 \times -1) = -5$

$(0 \times 1) + (1 \times 0) + (2 \times -1) + (5 \times 1) + (8 \times 0) + (9 \times -1) + (7 \times 1) + (2 \times 0) + (5 \times -1) = -4$

\* (Image - Vertical Edge detector)

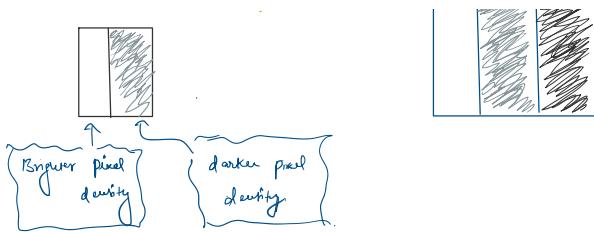
- tf. nn. Conv2d.

10	-	10	0	-	0
1	1	1	1	1	1
1	1	1	1	1	1
1	1	1	1	1	1
10	-	10	0	-	0

\*

1	0	-1
1	0	-1
1	0	-1





$$\textcircled{1} \quad \text{pixel value} \propto \frac{\text{bright pixels}}{\text{dark pixels.}}$$

### Horizontal Edge Detection

1	1	1
0	0	0
-1	-1	-1

More filters

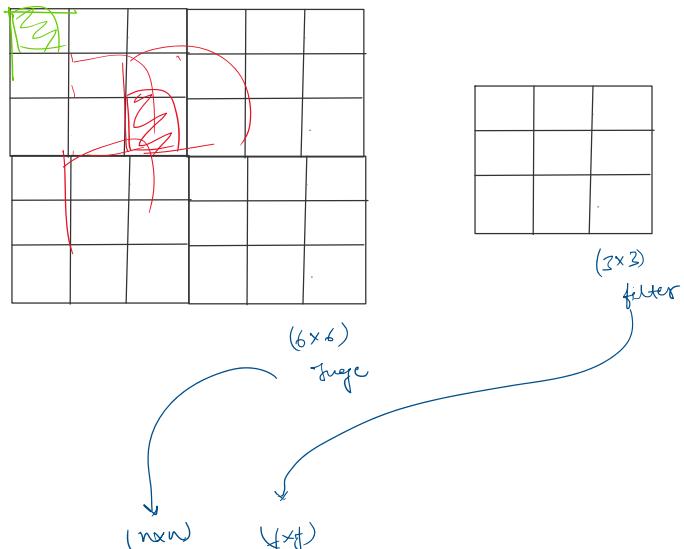
1	0	-1
2	0	-2
1	0	-1

3	0	-3
10	0	-10
3	0	-3

\* Sobel filter    \* Scharr filter.

Vertical Edge detector.

### Padding



$$(n-f+1) * (n-f+1)$$

$$(6,6) \xrightarrow[\text{f(3,3)}]{\text{Shrink}} (4,4) \xrightarrow{\dots} (1,1)$$

- Two ways shrinking.
- Top left pixel [0,0] is used only once but centre pixel is overlapped many times.

- Top left pixel  $[0,0]$  is used only once but centre pixel is overlapped many times.
- We are losing information from corner pixels.

Padding:

$$(6 \times 6) \iff (8 \times 8)$$

$(8 - 3 + 1) = 6 \times 6$  image (process image)

- padding of 0 :-
- padding = 1 :-  $\frac{(n-f+1)}{n+2p-f+1}$

### Valid and Same Convolutions

"Valid" :- No padding.

$$(n \times n) \rightarrow (f \times f) \rightarrow (n-f+1) \times (n-f+1)$$

"Same" :- Pad so that input size is same as output size.

$$(n+2p-f+1)$$

$$n+2p-f+1 = \cancel{x}$$

$$p = (f-1)/2$$

\*  $f$  is usually odd.

$f \in 2n+1$

### Strided Convolution

Stride  $\iff$  Step.

Stride = 2 (say).

Step key (Stride)

$$\text{padding: } p \quad \text{stride: } s \quad \left[ \frac{n+2p-f}{s} \right] + 1 \times \left[ \dots \right]$$

image:  $n$ .  
if this is not integer then floor ( $\lceil \rceil$ )

### Cross-Correlation v/s Convolution

- ① Sometimes we flip the filter (vertically and horizontally) and this process is called "Convolution".  
and process we perform is called cross-correlation.

$$f \rightarrow \begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 2 \\ -1 & 9 & 7 \end{bmatrix}$$

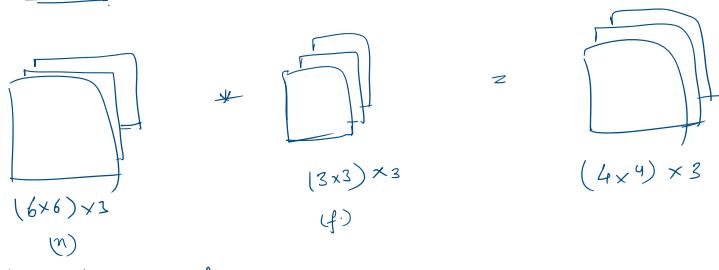
$$f^* \rightarrow \begin{bmatrix} 7 & 9 & -1 \\ 2 & 0 & 1 \\ 5 & 4 & 3 \end{bmatrix}$$

### Convolution Over Volumes

$\downarrow \leftarrow \downarrow \rightarrow \curvearrowright$

## Convolution Over Volumes.

### ① On RGB



height  $\times$  width  $\times$  channels



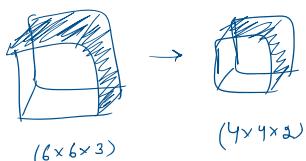
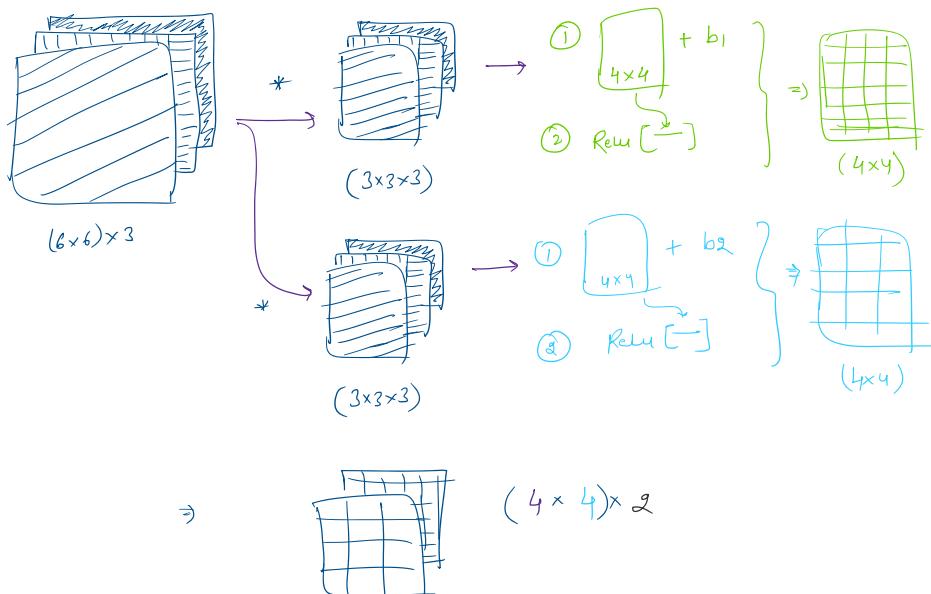
\* Channels of image = Channels on filters.

$$(n \times n \times n_c) * (f \times f \times n_c) \Rightarrow (n-f+1) \times (n-f+1) \times n_c$$

$$(6 \times 6) \times 3 * (3 \times 3) \times 3. \quad 4 \times 4 \times 2$$

### ② Assuming stride=1, padding=None.

## One Layer of CNN.



1 Layer of ConvNet.

- If layer  $l$  is a convolutional layer  $\Rightarrow$

$f^{[l]}$  = filter size.

$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_c^{[l]}$  = No. of filters.

Input :  $(n_h^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]})$

Output :  $(n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]})$

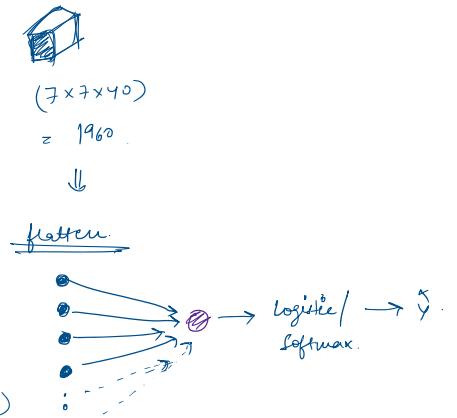
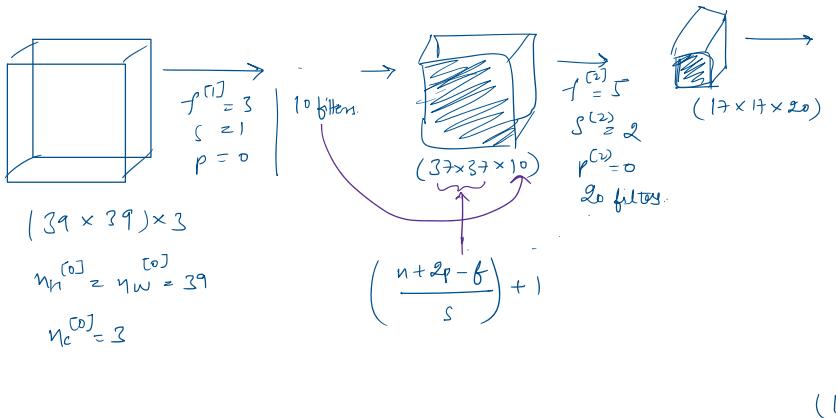
$$n^{[l]}_{(h/w)} = \left\lceil \frac{n^{[l-1]}_{(h/w)} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rceil$$

$$\begin{aligned}
 f^{(c)} &= \text{filter size.} \\
 p^{(c)} &= \text{padding} \\
 s^{(c)} &= \text{stride} \\
 n_c^{(c)} &= \text{No. of filters.} \\
 \text{Size of filter} &\Rightarrow f^{(c)} \times f^{(c)} \times n_c^{(c)}
 \end{aligned}$$

$$\begin{aligned}
 \text{Input: } & (n_H^{(c)} \times n_W^{(c)} \times n_C^{(c)}) \\
 \text{O/p: } & (n_H^{(c)} \times n_W^{(c)} \times n_C^{(c)}) \\
 n_H^{(c)} = & \left\lceil \frac{n_H^{(c)} + 2p^{(c)} - f^{(c)}}{s^{(c)}} + 1 \right\rceil
 \end{aligned}$$

$$\begin{aligned}
 \underline{\text{Activations}} &\Rightarrow a^{(c)} \rightarrow n_H^{(c)} \times n_W^{(c)} \times n_C^{(c)} \\
 A^{(c)} &\rightarrow m \times a^{(c)} \\
 \underline{\text{Weights}} &\Rightarrow (f^{(c)} \times f^{(c)} \times n_C^{(c)}) \times n_C^{(c)} \\
 &\quad \underbrace{(f^{(c)} \times f^{(c)} \times n_C^{(c)})}_{1 \text{ filter}} \quad \downarrow \quad (\# \text{filters}) \\
 \underline{\text{bias}} &\Rightarrow n_C^{(c)} \quad (1 \times 1 \times 1 \times n_C^{(c)})
 \end{aligned}$$

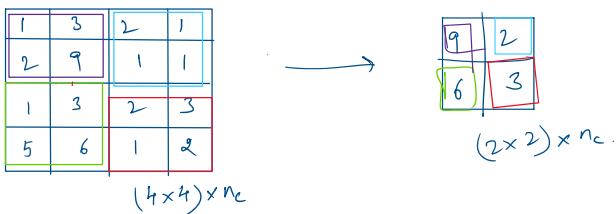
### Simpler Convolution N/w Example.



### Types of N/w in CNN:

- Convolution. (CONV)
- Pooling Layer. (POOL)
- Fully Connected (FC)

### Pooling layers. : Max pooling



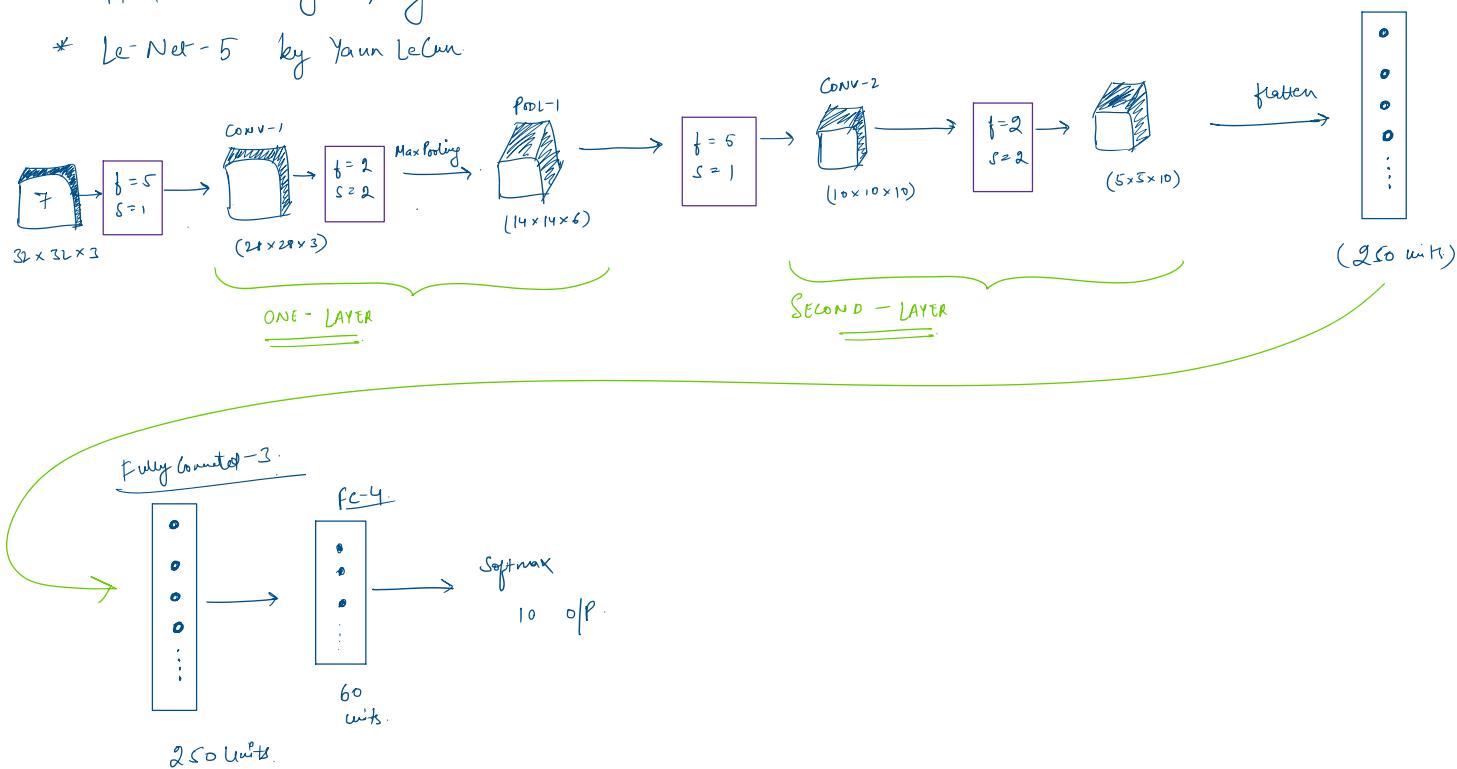
- Hyperparameters  $\Rightarrow$
- \* filter,  $f = 2$
  - \* stride,  $s = 2$
  - \* Take max over  $2 \times 2$  regions.
  - \* Takes no parameters.

$$\left( \frac{n+2p-f}{s} \right) + 1$$

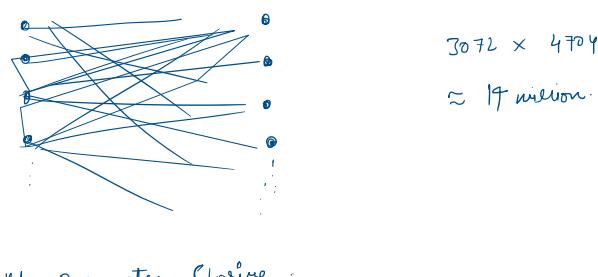
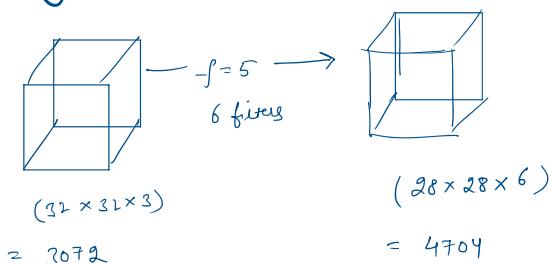
- Less used : Average pooling.

### CNN Example

- Handwritten digit recognizer.
- \* LeNet-5 by Yann LeCun.



### Why Convolutions? ?



With n - to - 1 connections.

## #1 Parameter Sharing :-

A feature detector that's useful in one part of the image is probably useful in another part of image.

## #2 Sparcity of Connections :-

In each layer, each o/p value depends only on a small number of units.

1 o/p unit is connected with  $\frac{6}{36}$  units.  
 $\downarrow$   
 $(^n)$      $(n \times n)$

### Classic N/w's.

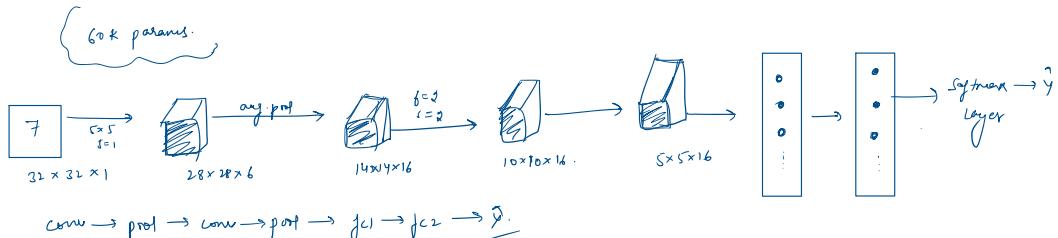
- LeNet-5
- AlexNet
- VGG-16

→ ResNet (Residual N/w).

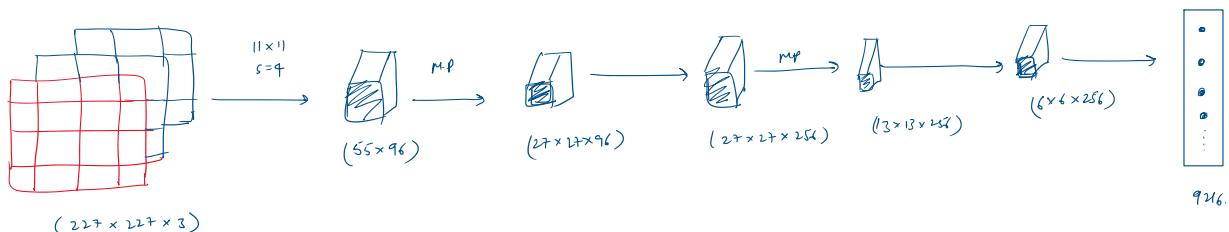
152 layer n.

→ Inception.

### Le Net - 5



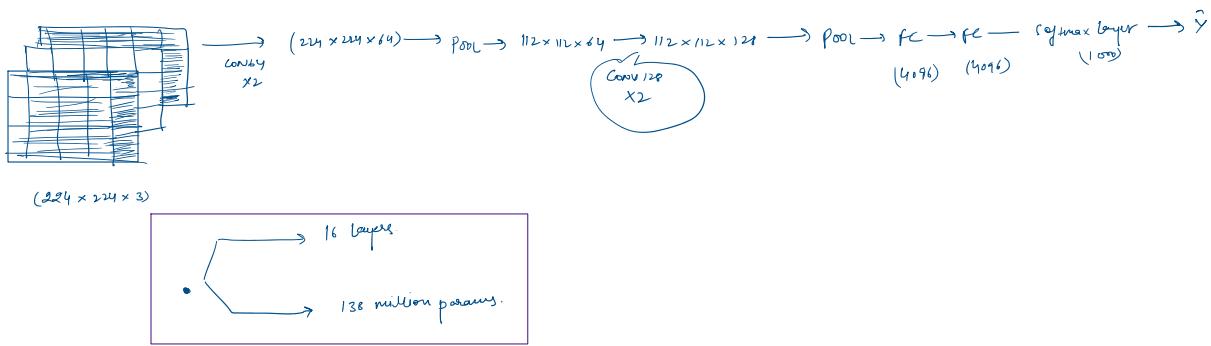
### AlexNet



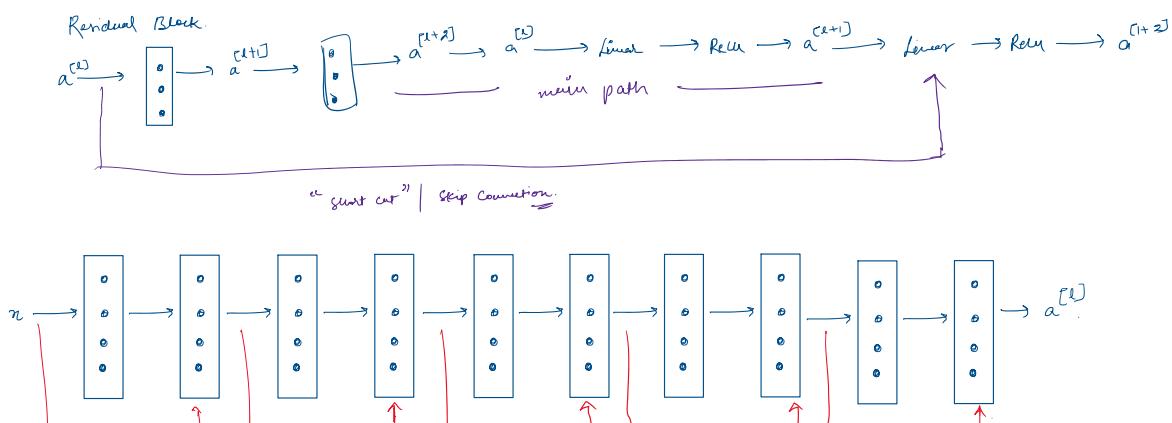
### VGG-16

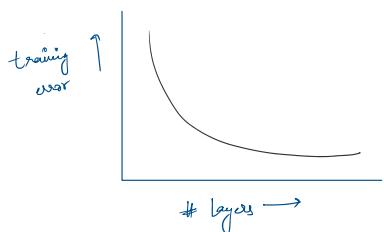
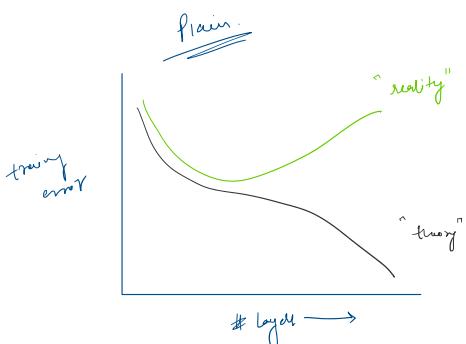
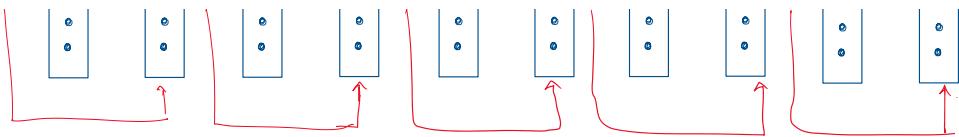
CONV =  $3 \times 3$  filter,  $s=1$ , same.

MAX POOL =  $2 \times 2$ ,  $s=2$ .

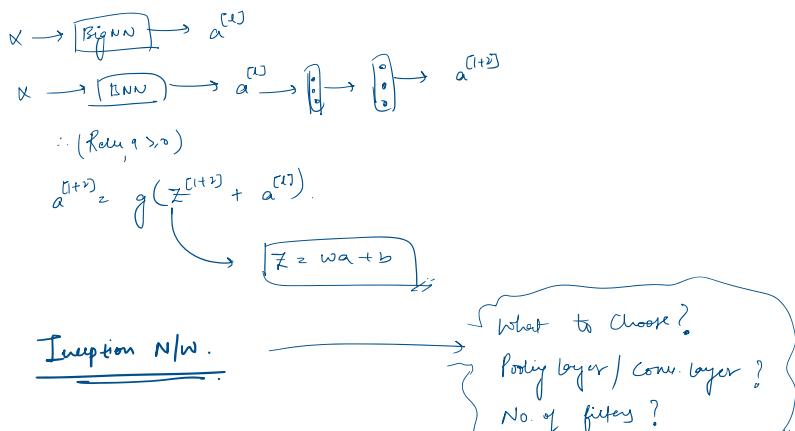


### ResNets - Residual N/w's

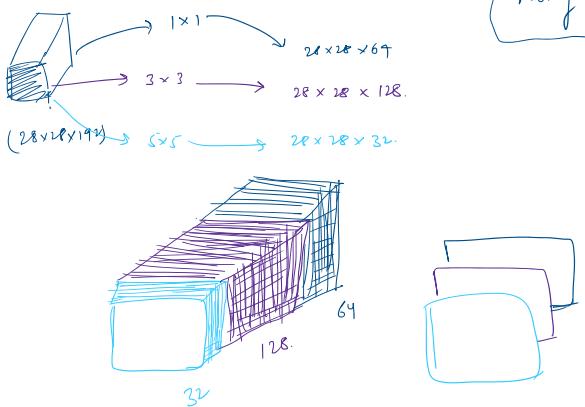




### Why ResNets Work?



### Inception N/W.



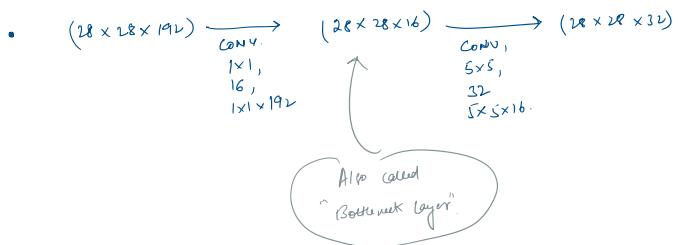
### Computational Cost

$$(28 \times 28 \times 192) \xrightarrow{\text{CONV } 5 \times 5, \text{ same, } 32} (28 \times 28 \times 32)$$

$$(28 \times 28 \times 192) \times (5 \times 5 \times 192)$$

$\approx 120$  million.

Reduce cost by  $1 \times 1$  convolutions!



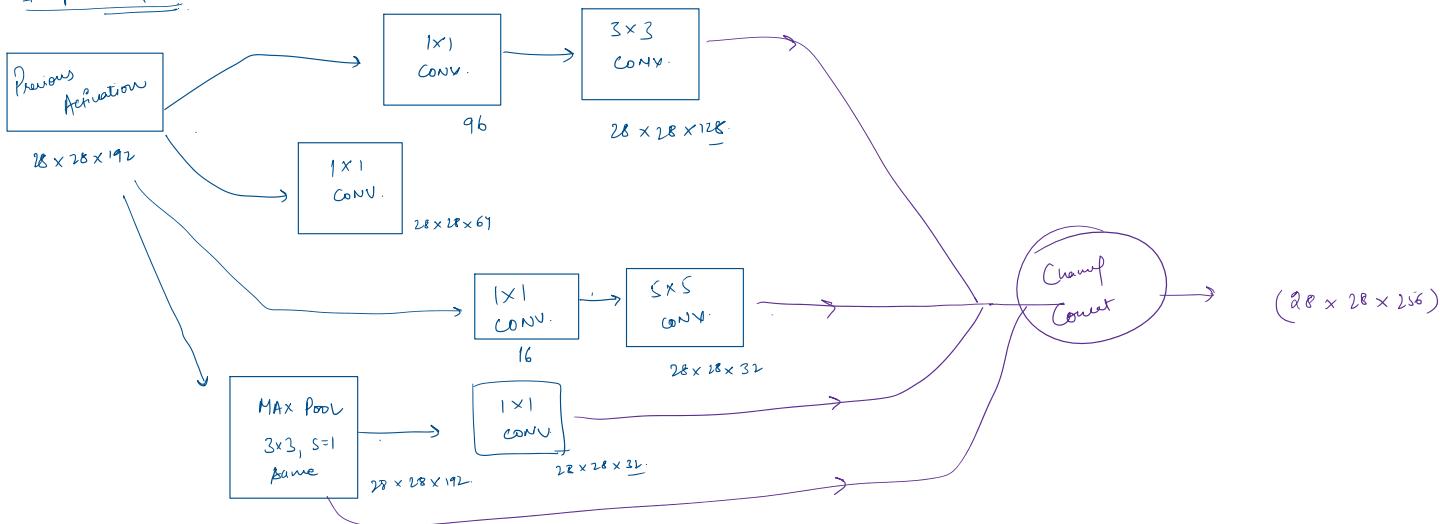
Convolutional Cost.  $\Rightarrow (28 \times 28 \times 16) \times 192 \approx 2.4$  million

+

$$(28 \times 28 \times 256) \times (5 \times 5 \times 16) \approx 10$$
 million

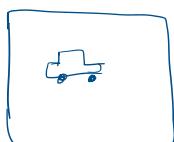
$\approx 12.4$  million

### Inception Net



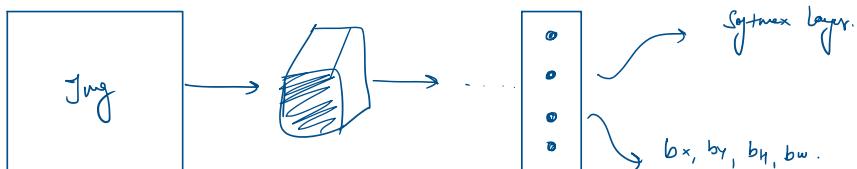
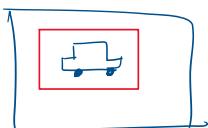
## Object Localization

Image classification



1

Image classification with localization



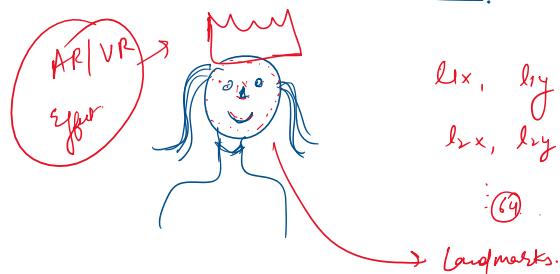
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Probability it is an object.

boundary boxes.

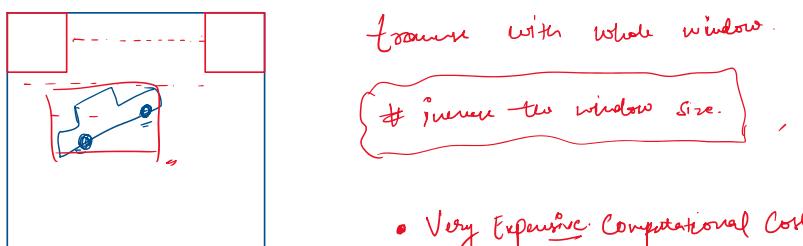
Class Names.

## Landmark Detection



## Object Detection

Sliding window detection

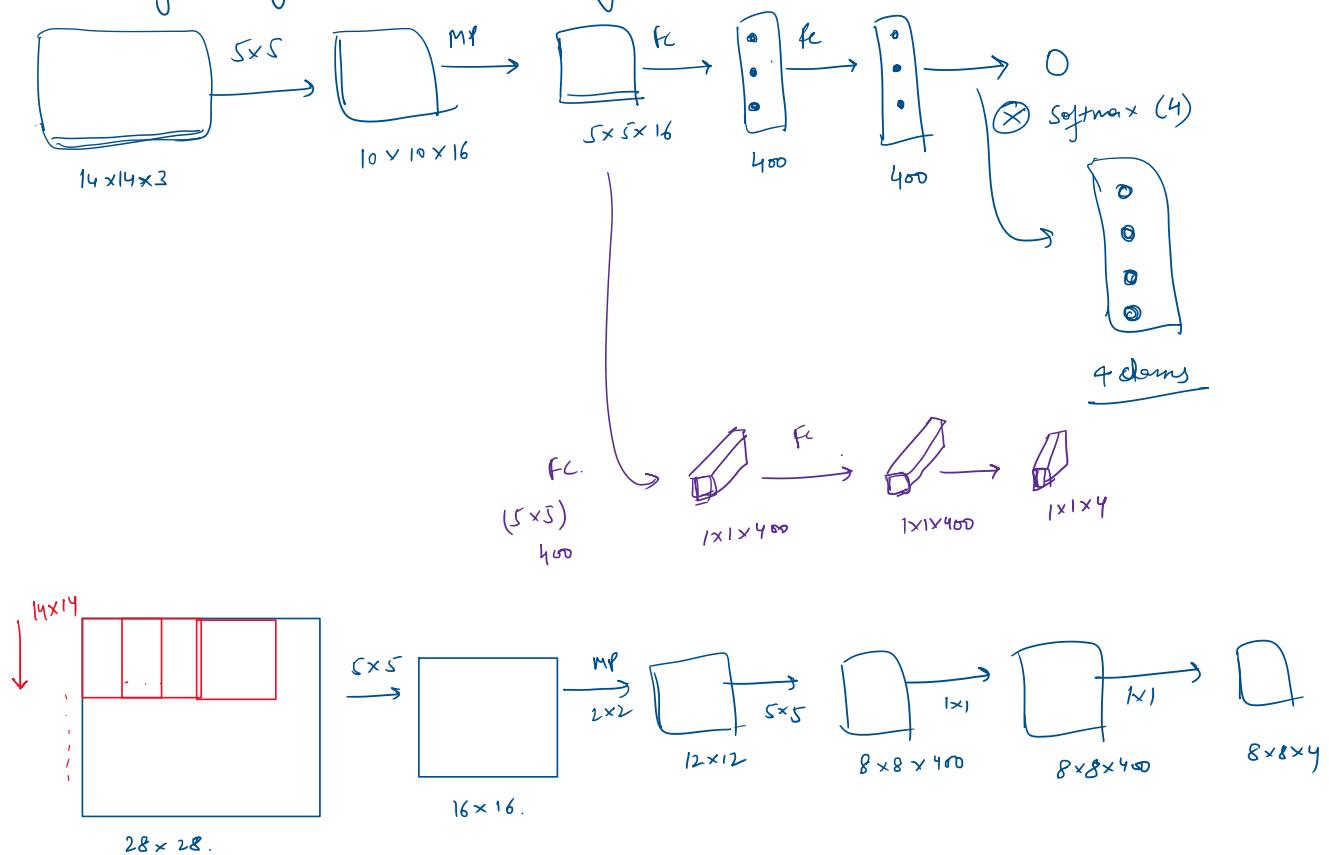




- Very Expensive Computational Cost

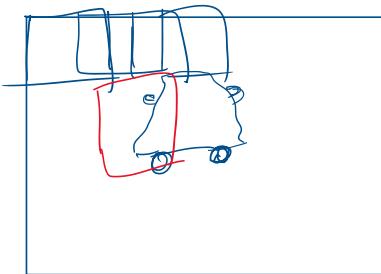
## Convolution implementation of Sliding Window.

- Turning FC layers into Convolutional layers.



## Bounding box prediction.

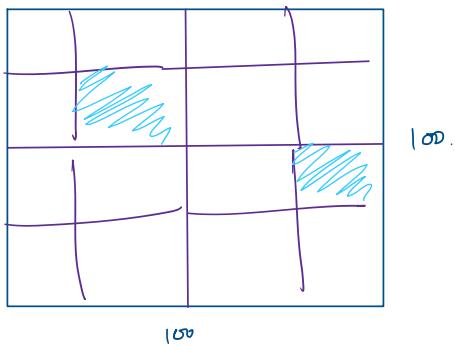
- Output accurate bounding boxes.



\* None of the bounding box is Accurate!

## YOLO Algorithm.

You only look once



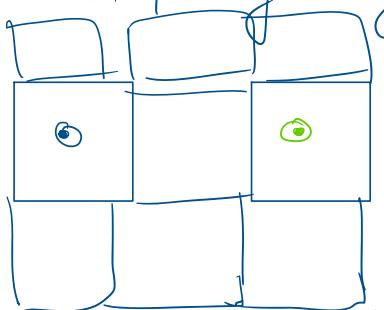
- \* divide into grids.
- ~ Object to be detected.

Labels for training

for each grid cell -

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

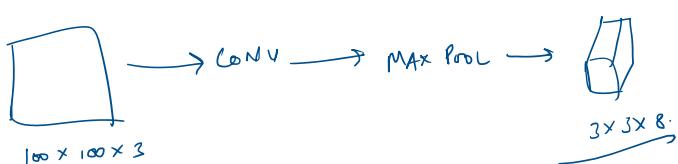
① Mid point of each grid cell is assigned if object is detected.



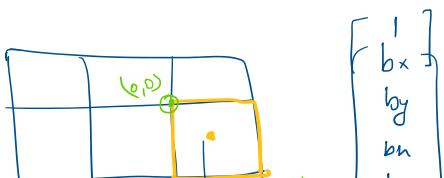
Total Volume of O/P  $\Rightarrow$

$$3 \times 3 \times 8$$

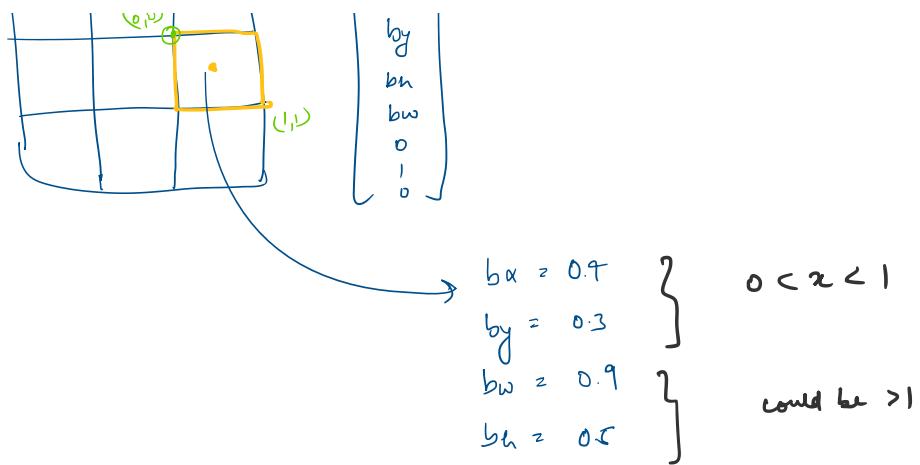
8 dimensional  
o/p vector.



Specifying the bounding box.

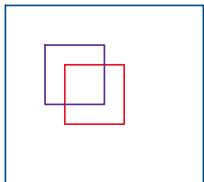


$$\begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \end{bmatrix}$$



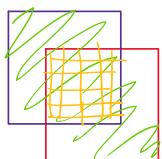
## INTERSECTION OVER UNION.

### Evaluating Object Localization



- ◻  $\rightarrow$  ground truth bounding box.
- ◻  $\rightarrow$  Algorithm Output.

IoU



# Union

# Intersection

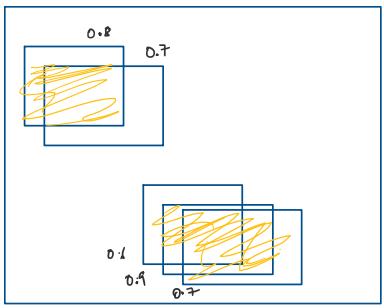
$$\frac{\text{Size of intersection} \#}{\text{Size of union} \#}$$

*Correct if IoU  $\geq 0.6$*

- (\*) IoU is the measure of Overlap.

## NON-MAX SUPPRESSION.

- (\*) Algorithm can detect same object multiple times, NMS makes sure that object is detected only once.

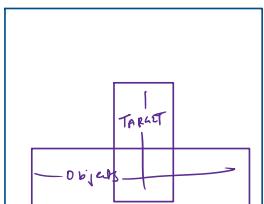


" $p(c)$  = Probability of Object."

- Discard all boxes with  $p(c) \leq 0.6$   
With the remaining boxes  $\Rightarrow$
- Pick the boxes with maximum  $p(c)$ .
- Discard any remaining box with  $IoU > 0.5$  with the box  
Output in the previous step.

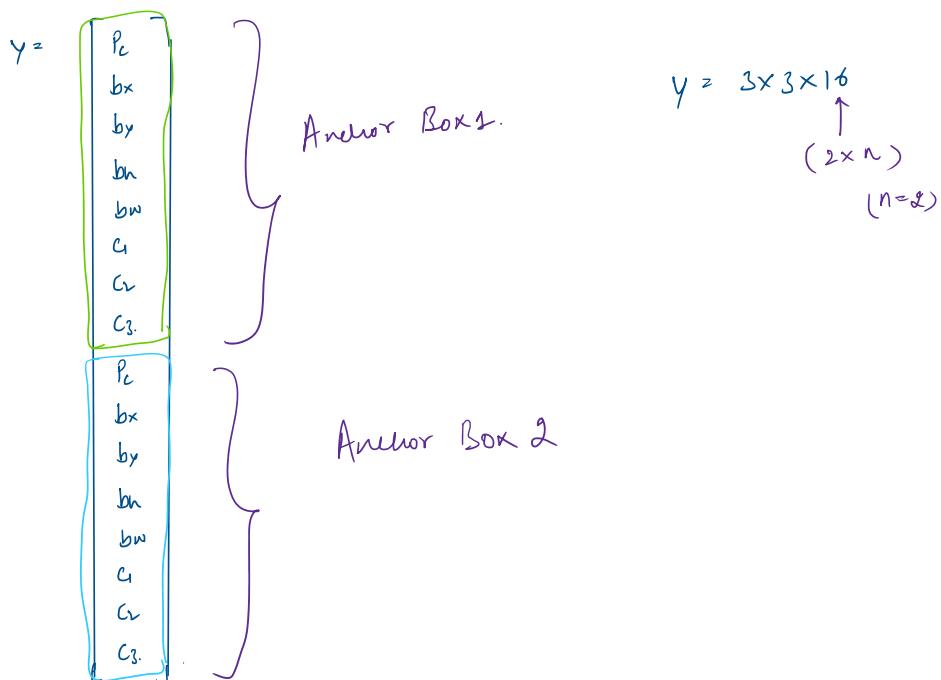
## ANCHOR BOXES.

(X) Only one object in the gridcell can be identified



Predefine two ( $n$ ) Anchor Boxes

- Anchor Box 1
- Anchor Box 2

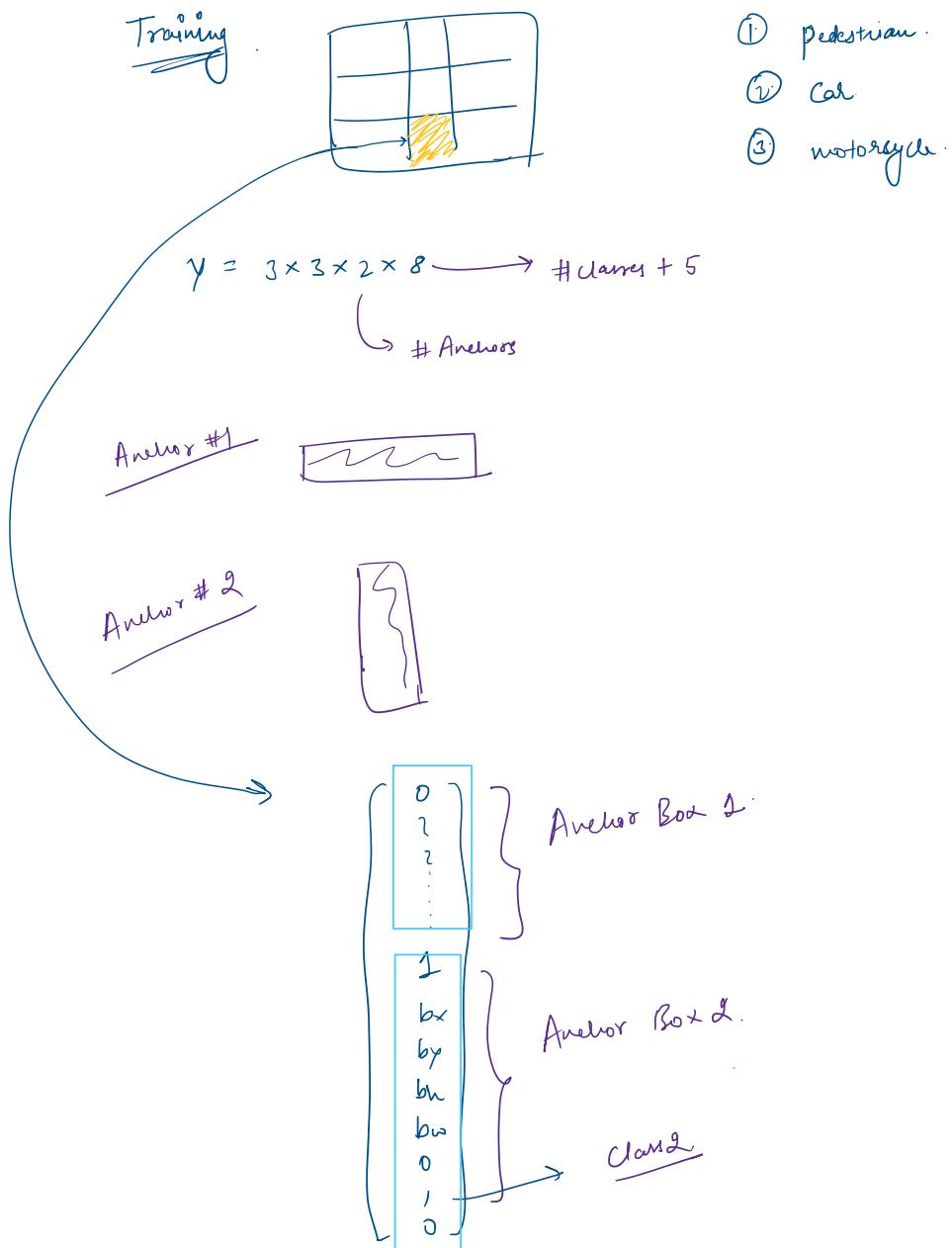


## YOLO ALGORITHM : Putting it together.

Training



① Pedestrian



Outputting the non-max Suppnd O/P :

- For each grid cell, get two predicted bounding boxes.
  - Get rid of low probability predictions.
  - For each class, use non-max suppression to generate final predictions.



## FACE RECOGNITION.

### ① Verification

→ Input image.

→ O/p whether the image is of claimed person (0/1).

### ② Recognition.

→ Has database of  $K$  persons.

→ get an input image.

→ O/p Id of that person if in db.

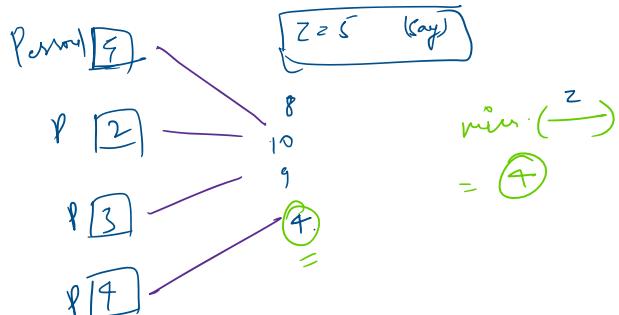
## ONE SHOT LEARNING

"Learn by only single input image".

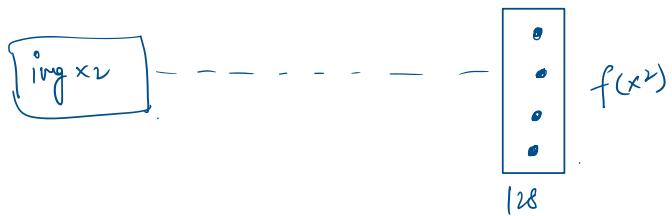
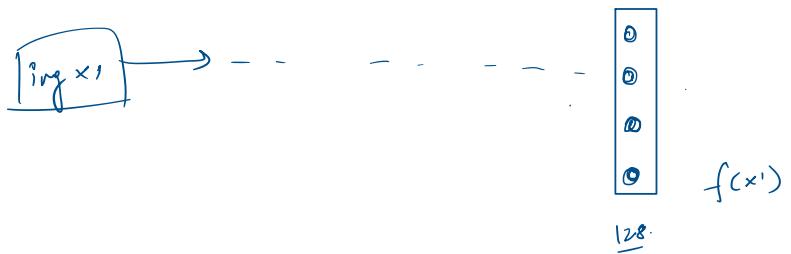
\* Learn similarity fx.

$d(\text{img}_1, \text{img}_2) = \text{degree of difference b/w images img}_1 \text{ and img}_2$ .

if  $\begin{cases} d \leq z & \text{"Same"} \\ d > z & \text{"different"} \end{cases}$   $\quad \textcircled{z} = \text{Tau (some hyperparameter)}$



## SIAMESE NETWORK.



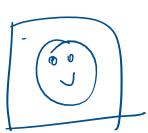
$$d(x_1, x_2) = \|f(x_1) - f(x_2)\|_2^2$$

=

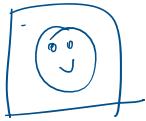
Comparing two images and calculating d is sometimes called  
"Siamese Network".

"Deep Face"

### Triplet Loss fn.

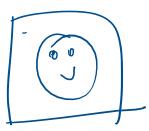


Anomalous image.



Positive image.

$[z \ll]$



Anomalous image.



Negative image.

$[z \gg]$

Triplet : { Anchor , A  
+ve , P  
-ve , N. }

- $\| f(A) - f(P) \|^2 \leq \| f(A) - f(N) \|^2$

$d(A, P)$                              $d(A, N)$

$$d(A, P) - d(A, N) \leq 0$$

$$\| f(A) - f(P) \|^2 - \| f(A) - f(N) \|^2 \leq 0 - \alpha$$

or

(hyperparameter)

$$\left[ \| f(A) - f(P) \|^2 - \| f(A) - f(N) \|^2 \right] + \alpha \leq 0$$

(Margin)

$\alpha = 0.2 \text{ (say)}$

Loss fx.

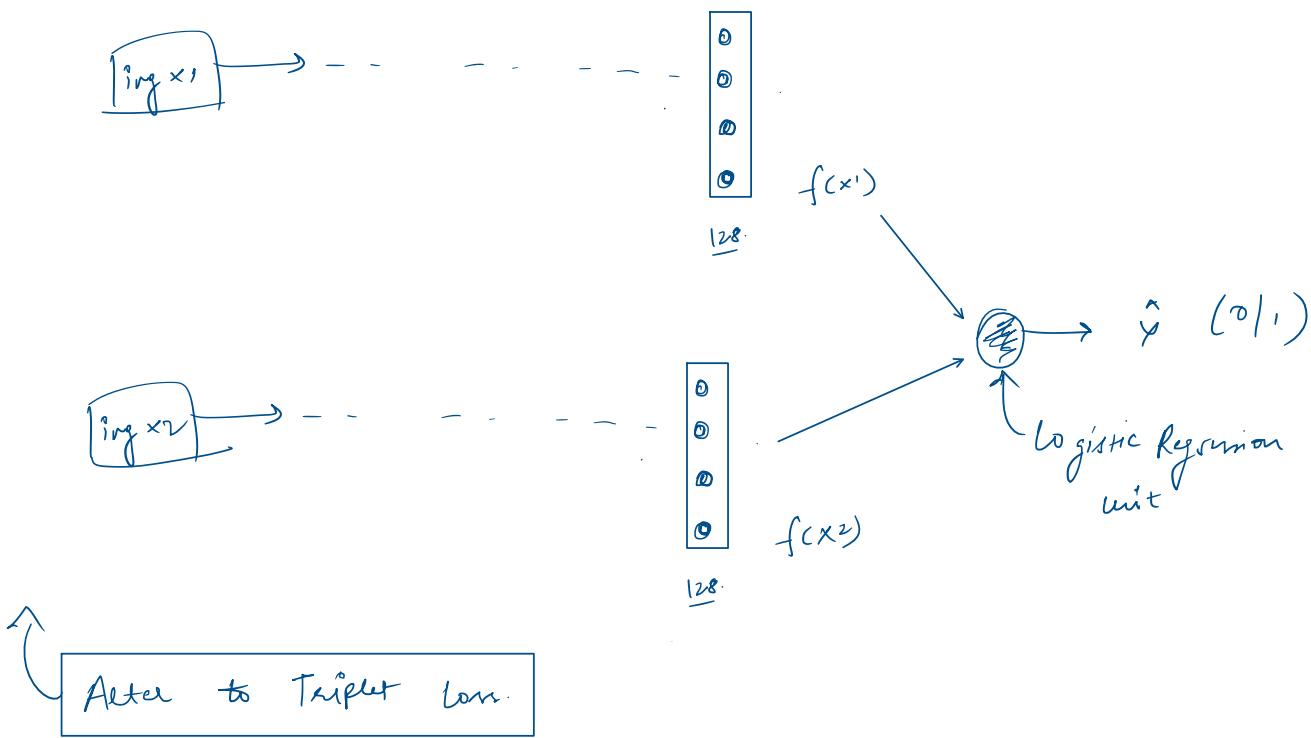
- given three images  $A, P, N$ .

$$\ell(A, P, N) = \max \left( \| f(A) - f(P) \|^2 - \| f(A) - f(N) \|^2 + \alpha, 0 \right)$$

$\leq 0$

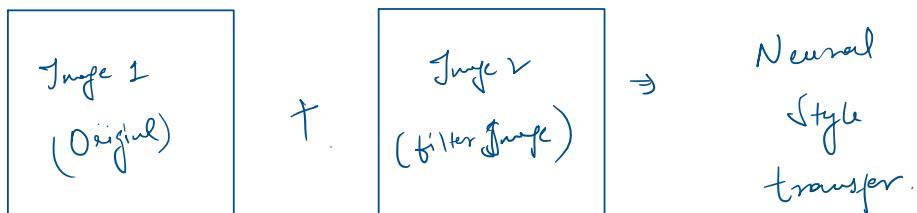
$$J = \sum_{i=1}^m \ell(A_i, P_i, N_i)$$

Face Verification & binary Classification.



$$\hat{y} = \sigma \left( \sum_{k=1}^{128} w_k |f(x^i)_k - f(x^j)_k| + b \right)$$

## NEURAL STYLE TRANSFER

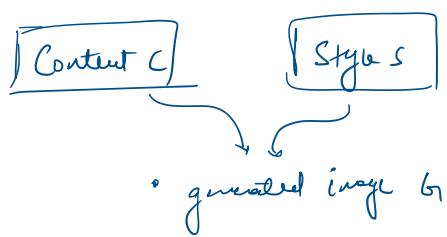


### ④ deep Convnets learning



Pick a unit in layer 1. Find the 9 image patches that maximize the units activation. Repeat for other units.

## Cost fx.



$$J(G) = \alpha J_{\text{content}}(C, G)$$

$\downarrow$   
similarity

+

$$\beta J_{\text{style}}(S, G)$$

Initiate  $G$  randomly,

$$G_1 : 100 \times 100 \times 3$$

$$G_t := G_{t-1} - \frac{\partial}{\partial G_t} J(G_t)$$

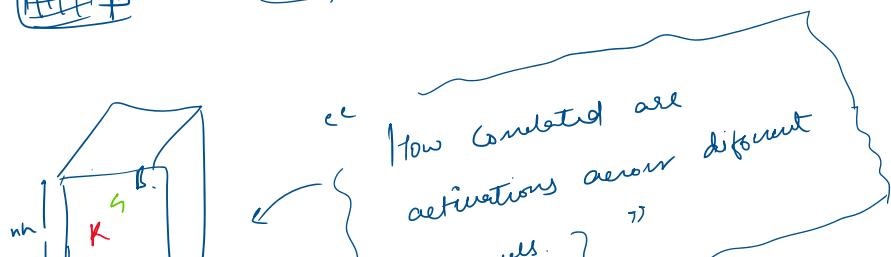
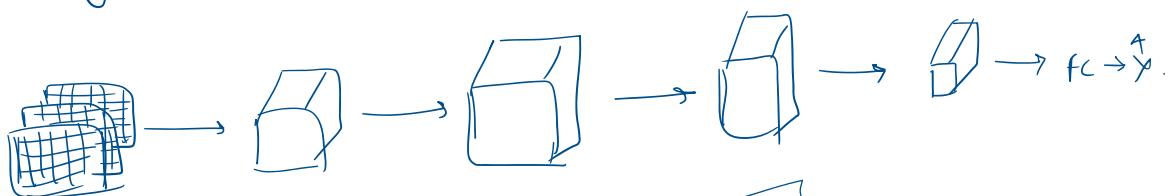
## Content Cost fx.

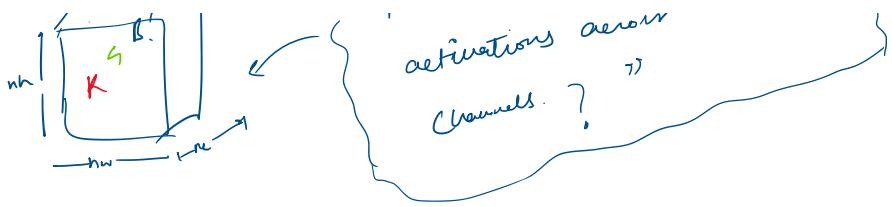
$$J(G) = \underbrace{\alpha J_{\text{content}}(C, G)} + \beta J_{\text{style}}(S, G)$$

- Say you use layer  $L$ .
- Use pre-trained ConvNet (say VGG).
- Let  $a^{[L](C)}$  and  $a^{[L](G)}$  be activation on layer  $L$  on the image.
- If  $a^{[L](C)} \approx a^{[L](G)}$  are similar, both images have similar content.

$$J_{\text{content}}(C, G) = \frac{1}{2} \| a^{[L](C)} - a^{[L](G)} \|_2^2$$

## Style Cost fx.





### Style Matrix

Let  $a_{ijk}^{(e)} = \text{activation at } (i, j, k)$

$G^{(e)}$  is  $n_e \times n_c^{(e)}$ .

$$G_{KK1}^{(L)(G)} = \sum_{i=1}^{nh} \sum_{j=1}^{nw} \begin{pmatrix} a_{ijk}^{(L)(G)} & a_{ijk}^{(L)(G)} \end{pmatrix}.$$

(Correlation)

$$G_{KK1}^{(L)(S)} = \sum_{i=1}^{nh} \sum_{j=1}^{nw} \begin{pmatrix} a_{ijk}^{(L)(S)} & a_{ijk}^{(L)(S)} \end{pmatrix}.$$

$G$  = generated image.

$S$  = style image

"gram matrix".

$$\boxed{J_{\text{Style}}(S, G) = \| G^{(L)(S)} - G^{(L)(G)} \|_F^2}$$

$$= \frac{1}{(2n_H^{(e)} n_W^{(e)} n_C^{(e)})^2} \sum_k \sum_{kl} (G_{kk1}^{(L)(S)} - G_{kk1}^{(L)(G)})$$

$$\text{or } J(G) = \alpha [J_{\text{Content}}(L, G)] + \beta [J_{\text{Style}}(S, G)].$$

### 1d and 3d generalization

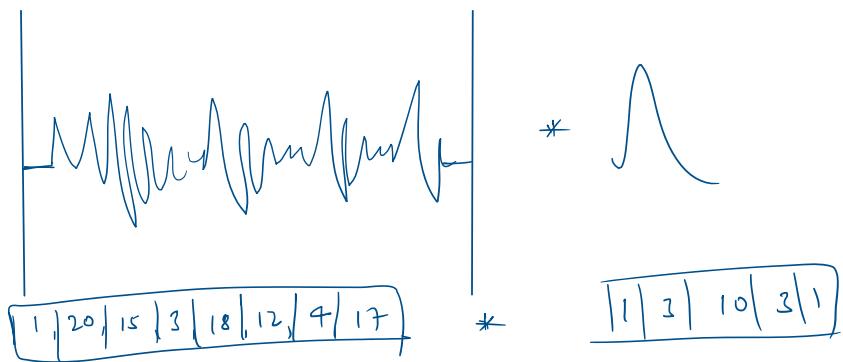
$$\begin{array}{ccc} \text{2d image} & * & \text{2d filter} \\ (14 \times 14) & & (5 \times 5) \end{array} = (10 \times 10)$$

$$\begin{array}{l} \text{2d image} \\ (14 \times 14) \end{array} * \begin{array}{l} \text{2d filter} \\ (5 \times 5) \end{array} = (10 \times 10)$$

# Same can be applied to 1d image.

ex. Ecg

electro cardio graphy



## Why Sequence Models?

- Examples of Sequence data :-

① Speech recognition :-  → "The quick brown fox jumps over the lazy dog".

② Music generation :-  $\phi$  → 

③ Sentiment classification :- "bad movie" → 1★

④ DNA Sequence Analysis → AGCCCTTACCG

⑤ Machine translation → "Hermosa" → "Beautified"

⑥ Video activity recognition →  → "Running".

⑦ Name entity recognition → "Yesterday Harry Potter met Hermoine granger".

## NOTATIONS

$x$ : "Harry Potter & Hermoine granger invented a new spell".  
 ↓      ]  
 Name entity recognition

$y$ : 1 1 0 1 1 0 0 0 0

- input is sequence of 9 words, 9 sets of features to represent these 9 words.

$x^{(t)}$ ,  $x^{(1)}$ ,  $x^{(2)}$  ---  $x^{(9)}$   $T_x = 9$   
 "temporal Sequence".

$y^{(t)}$ ,  $y^{(1)}$ ,  $y^{(2)}$  ---  $y^{(9)}$   $T_y = 9$

$x^{(i)(t)}$ ,  $T_x^{(i)}$

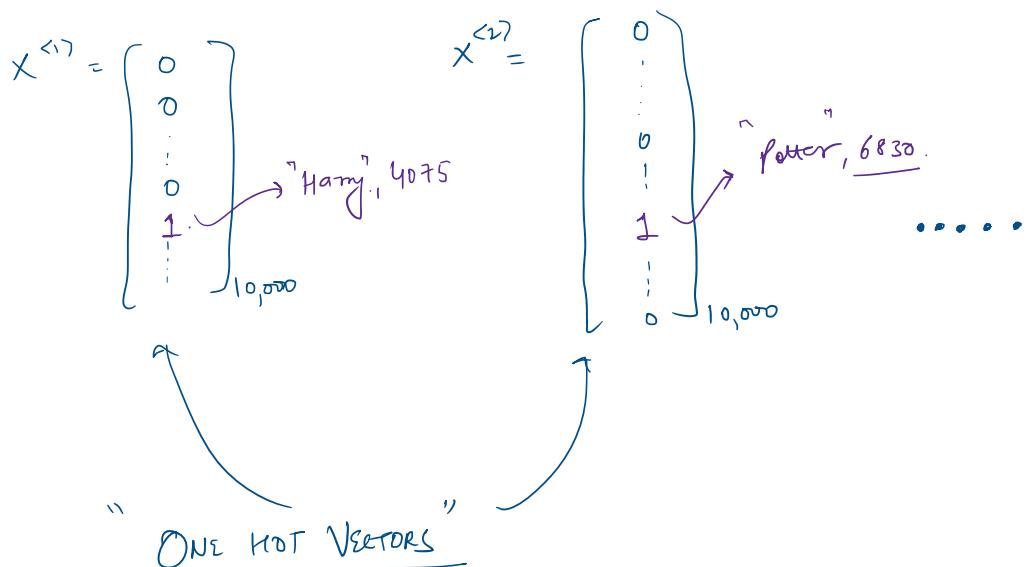
$$x^{(i)<\tau}, \quad T x^{(i)}$$

$$y^{(i)<\tau}, \quad T y^{(i)}.$$

Vocabulary	dictionary
a	1
aeron	2
and	267
Harry	4075
Potter	6830
Zuker	10,000

(Assuming).

### One Hot representations



### Recurrent Neural N/w MODEL.

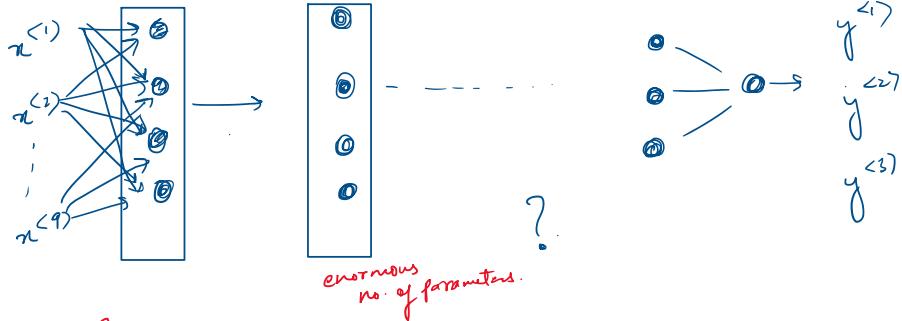
Why not a Standard Neural N/w ?

$x^{(i)}$   $\rightarrow$   $\boxed{\oplus}$

$\boxed{\oplus}$

$y^{(i)}$

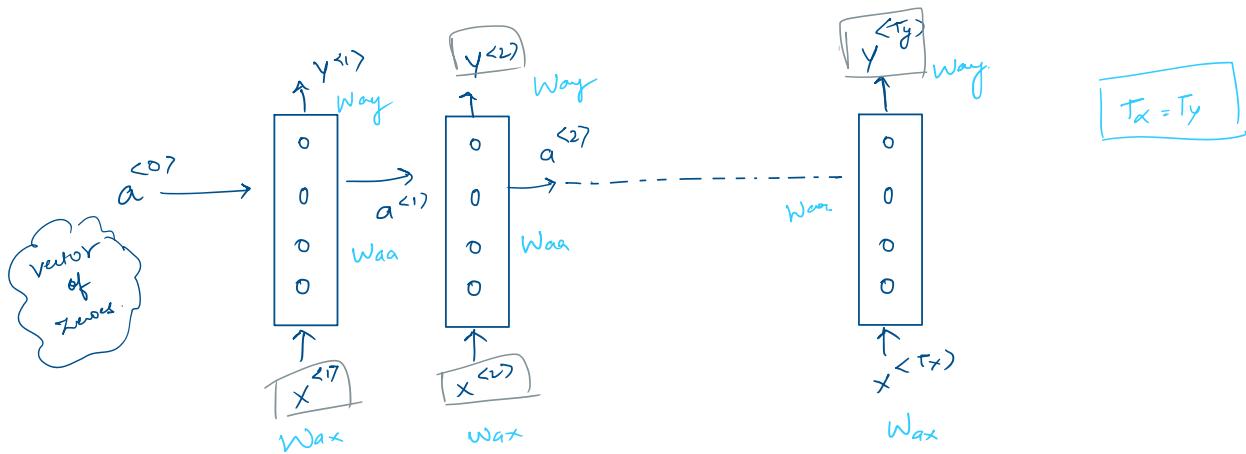
0 \_\_\_\_\_



### Problems!

- ①  $I/P$  and  $O/P$  can be different layers in diff example.
- ② doesn't share features learned across different portion of text.

### Recurrent Neural NW

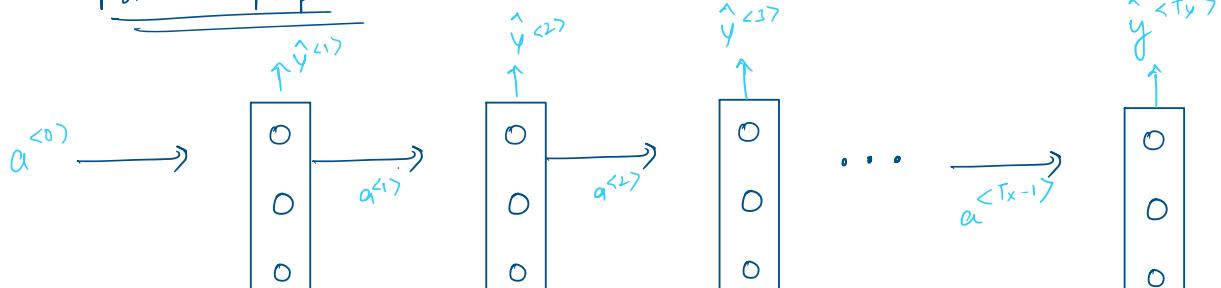


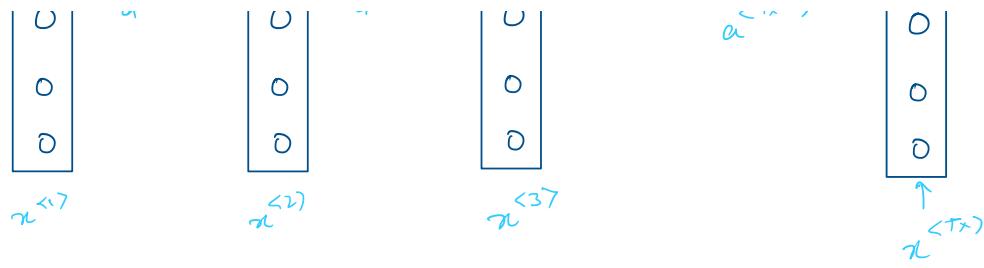
> While making predictions  $y^{(i)}$ , it also uses  $x^{(i)}$  and  $x^{(1)}$   
or  
( earlier in the sequence. )

ex: He said "Teddy Roosevelt was a great president." (✓)

He said "Teddy Bears are on Sale". (✗)

### Forward Prop.





$$a^{<0>} = \vec{0}.$$

$$a^{<1>} = g_1(W_{aa} \cdot a^{<0>} + W_{ax} \cdot x^{<1>} + b_a)$$

$$\hat{y}^{<1>} = g_2(W_{ya} \cdot a^{<1>} + b_y)$$

$\rightarrow \text{tanh}(\cdot) / \text{ReLU}$

$\rightarrow \text{Sigmoid}$ .

$$a^{<t>} = g(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya} \cdot a^{<t>} + b_y).$$

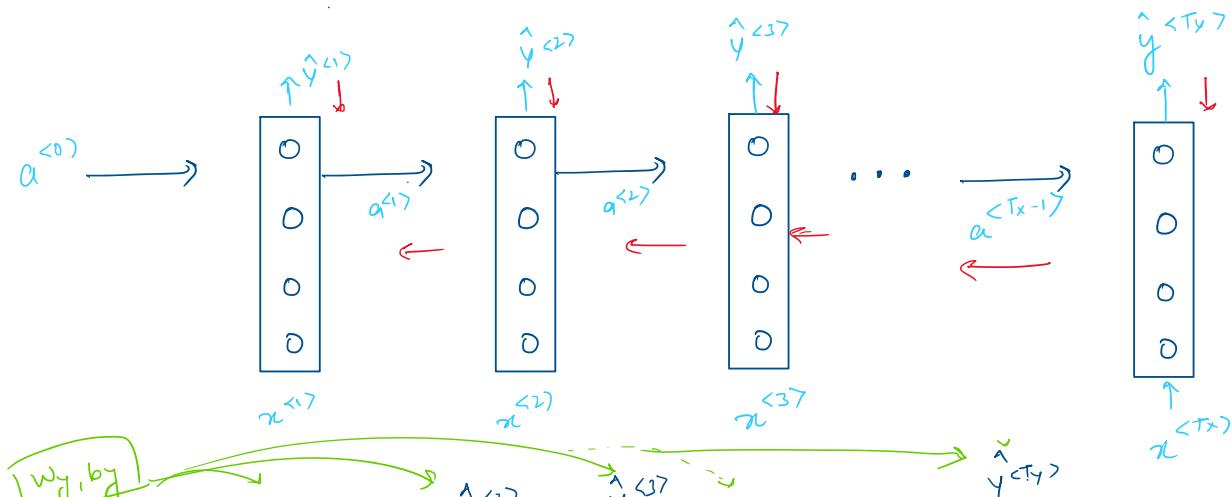
### Simplified RNN Notation

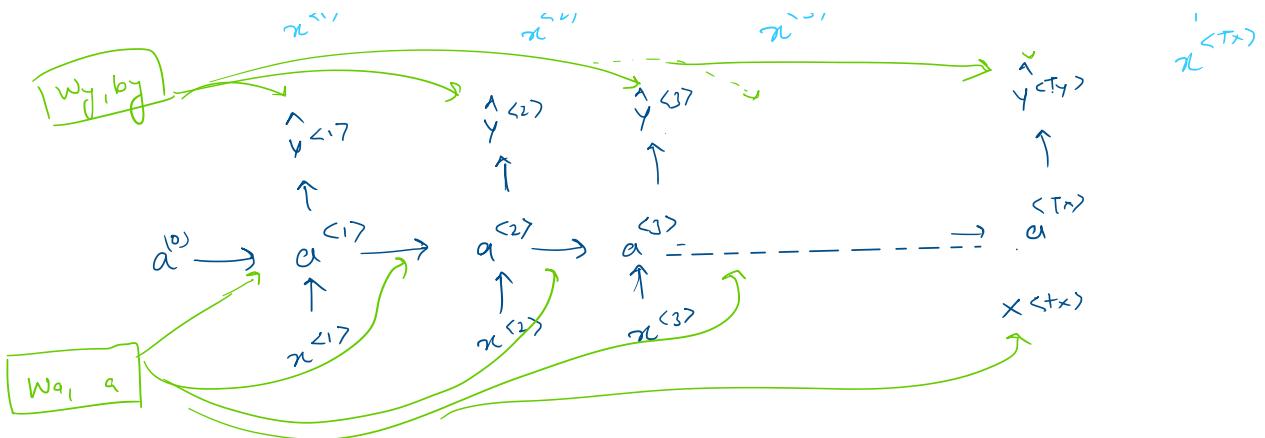
$$a^{<t>} = g(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a) \iff g(W_a(a^{<t-1>} + W_x x^{<t>})) + b_a$$

$$\hat{y}^{<t>} = g(W_{ya} \cdot a^{<t>} + b_y).$$

$$\begin{bmatrix} W_{aa} & | & W_{ax} \end{bmatrix} \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = \boxed{W_a(a^{<t-1>} + W_x x^{<t>})}$$

### Backprop. Through Time.



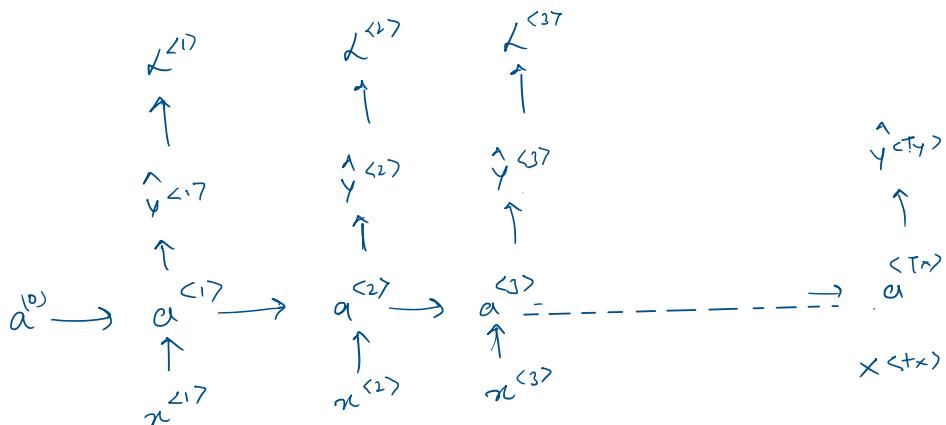


Loss fx.

$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1-y^{<t>}) \log (1-\hat{y}^{<t>})$$

"Cross entropy loss"

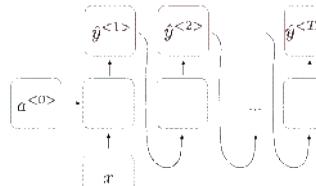
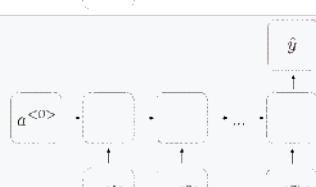
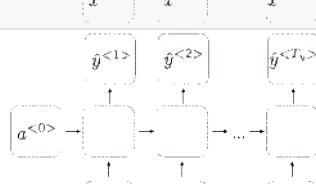
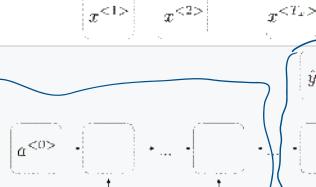
$$L(y, \hat{y}) = \sum_{t=1}^{Tx} L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$



Different Types Of RNN.

[when  $T_x = T_y$ ]

Type of RNN	Illustration	Example
One-to-one $T_x = T_y = 1$		Traditional neural network

One-to-many $T_x = 1$ , $T_y > 0$		Music generation
Many-to-one $T_x > 1$ , $T_y = 1$		Sentiment classification
Many-to-many $T_x = T_y$		Name entity recognition
Many-to-many $T_x \neq T_y$		Machine translation

From <<https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks#overview>>

Stanford University CS230

# Encoder

• 21 •

# LANGUAGE MODEL AND SEQUENCE GENERATION.

## Speech Recognition

"The apple and pear salad."

pern / pair. ?

$$\bullet \quad P(\text{The apple and pear School}) = 5.7 \times 10^{-10} \quad \checkmark$$

$$P(\text{The apple and pear scaled}) = 3.2 \times 10^{-13}$$

$$p(\text{future}) = ?$$

$$y^{(1)}, y^{(2)}, y^{(3)}, \dots \dots \dots y^{(Ty)}$$

## Language Building with RNN.

v v v

Training set  $\Rightarrow$  Large **Corpus** of english text.

Collection

Example Sentence :-

Cats average 15 hours of sleep a day.

① Tokenize (Create vocabulary/dictionary).

$y^{<1>} \quad y^{<2>} \quad y^{<3>} \dots \dots \dots \quad y^{<T>} \quad <\text{EOS}>$

append an extra token at the end of sequence  $<\text{EOS}>$   
end of sentence.

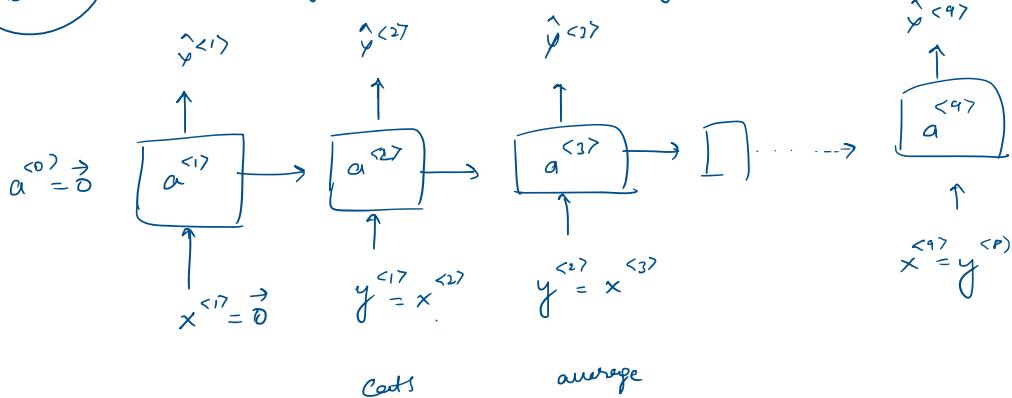
"The Egyptian Mau is a breed of cat."

if the word is not in vocabulary  
replace with  $\langle \text{UNK} \rangle$  token.  
 $\nearrow$  unknown word.

RNN Model

choose example

• Cats average 15 hours of sleep a day  $<\text{EOS}>$



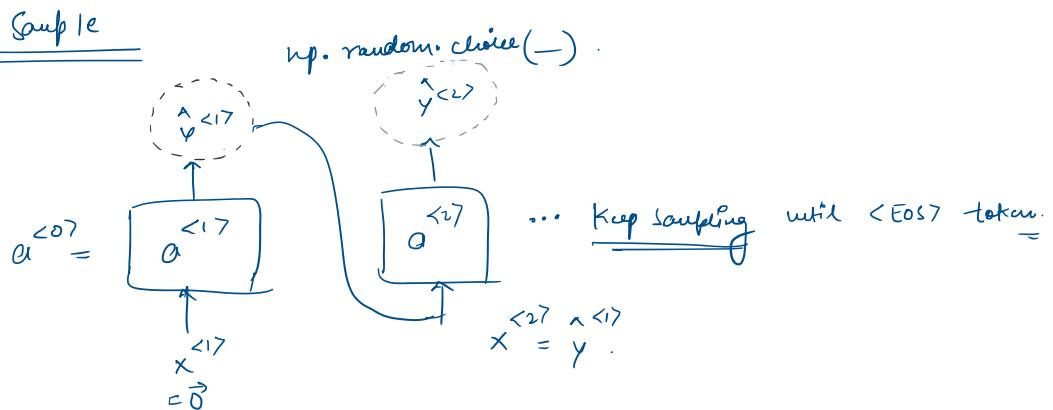
$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

$$\Rightarrow L = \mathcal{L}(\hat{y}^{<1>}, y^{<1>}) + \mathcal{L}(\hat{y}^{<2>}, y^{<2>}) + \dots + \mathcal{L}(\hat{y}^{<T>}, y^{<T>})$$

SAMPLING NOVEL SEQUENCES.

$$P(y^{<1>} | y^{<2>} | \dots | y^{<T>})$$

Sample

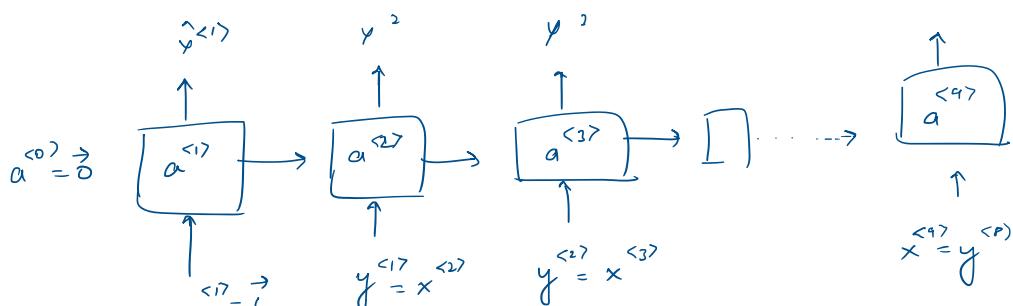


Character level words -

Vocabulary =  $[a, b, c, \dots, z, A-Z, ., , :, ;, " - \dots ]$

Very expensive.

## VANISHING GRADIENTS WITH RNN'S



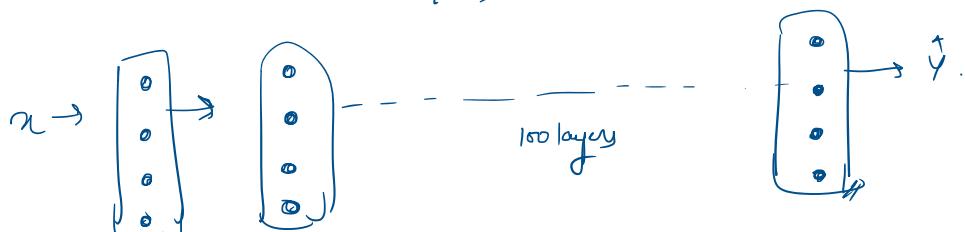
sentence : "The cat which already ate -----, was full."

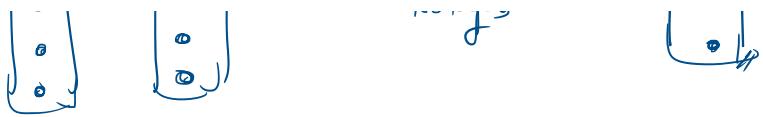
Cat is singular so

Cat : was

Cats : were.

ex. in deep neural net  $\Rightarrow$   
(100 layers)





forward prop.  $\xrightarrow{\quad}$  Backprop.

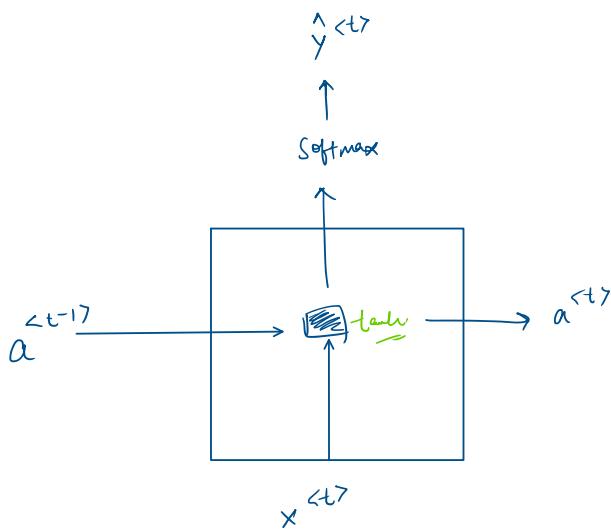
- if the layer is very deep, we will have hard time calculating gradient, affecting weights of earlier layers.

Vanishing / Exploding  
gradients  $\xrightarrow{\quad}$  by tuning weight.

$\rightarrow$  gradient clipping.

## GATED RECURRENT UNIT. (GRU)

$$a^{(t)} = \tilde{g}(W_a [a^{(t-1)}, x^{(t)}] + b_a)$$



"The cat which already ate . . . , was full."  
 {cat : was }  $\xrightarrow{\quad}$  {cats : were }

$a$  = memory cell.  $\left. \begin{array}{l} \\ \end{array} \right\}$  Provide memory to remember that the "cat" was singular/plural

$c$  = memory cell.       $\left\{ \begin{array}{l} \text{Provide memory to remember that the "cat" was singular/plural} \\ \text{at time } t \quad \text{for further consideration.} \end{array} \right.$

$$c^{<t>} \xrightarrow{o/p} a^{<t>}$$

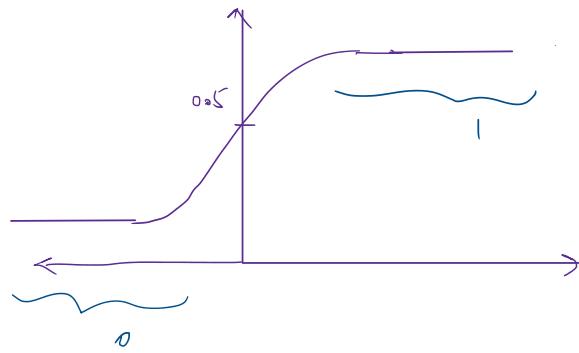
$$\tilde{c}^{<t>} = \tanh(w_c [c^{<t-1>} \times^{<t>}] + b_c).$$

Candidate for replacing  $c^{<t>}$ .

update gate

$$\gamma_u = (0, 1) = \sigma(w_u [c^{<t-1>} \times^{<t>}] + b_u).$$

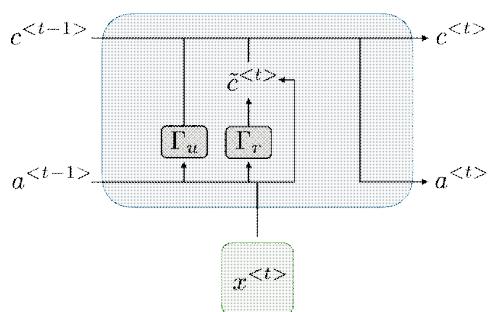
(gamma)



{  
singular : 1,  
plural : 0  
}

$$c^{<t>} = \gamma_u * \tilde{c}^{<t>} + (1 - \gamma_u) * c^{<t-1>} \quad \boxed{=}$$

FULL GRU



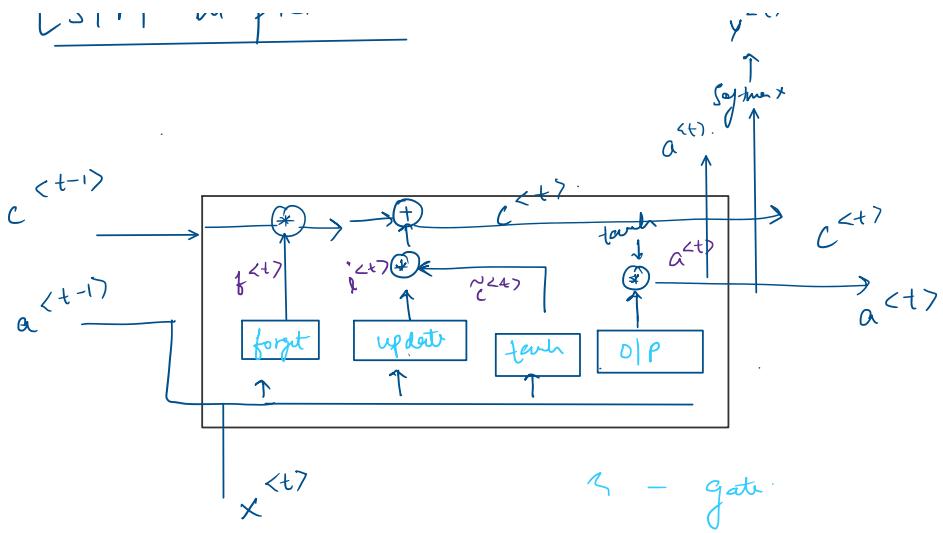
$$\begin{aligned}
 h & C^{t+1} = Y_u * \tilde{C}^{t+1} + (1 - Y_u) * C^{t+1} \\
 \tilde{h} & \tilde{C}^{t+1} = \tanh \left( w_c [Y_r * C^{t+1}, x^{t+1}] + b_c \right) \\
 r & Y_r = \sigma(w_r [C^{t+1}, x^{t+1}] + b_r) \\
 u & Y_u = (0, 1) = \sigma(w_u [C^{t+1}, x^{t+1}] + b_u) \\
 a^{t+1} & = C^{t+1}
 \end{aligned}$$

## LONG-SHORT TERM MEMORY. (LSTM)

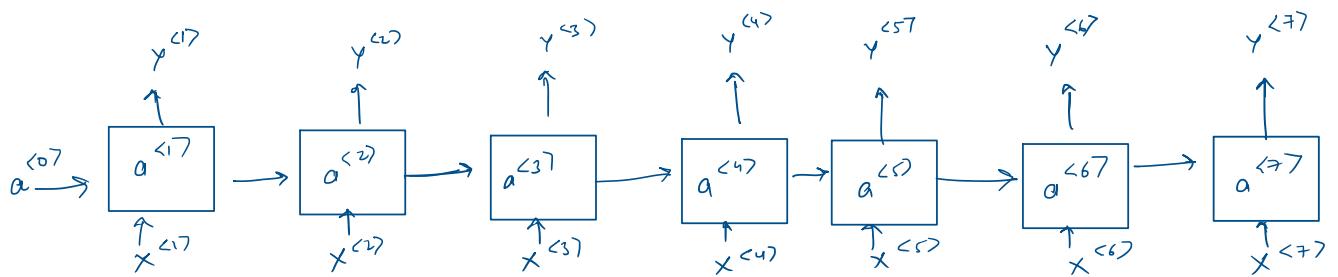
$$\begin{aligned}
 C^{t+1} &= Y_u * \tilde{C}^{t+1} + (1 - Y_u) * C^{t+1} \\
 \tilde{C}^{t+1} &= \tanh \left( w_c [a^{t+1}, x^{t+1}] + b_c \right) \\
 Y_r &= \sigma(w_r [C^{t+1}, x^{t+1}] + b_r) \\
 Y_u &= (0, 1) = \sigma(w_u [a^{t+1}, x^{t+1}] + b_u) \\
 Y_f &= \sigma(w_f [a^{t+1}, x^{t+1}] + b_f) \\
 V_o &= \sigma(w_o [a^{t+1}, x^{t+1}] + b_o) \\
 C^{t+1} &= [Y_u * \tilde{C}^{t+1}] + [Y_f * C^{t+1}] \\
 a^{t+1} &= Y_o * \tanh(C^{t+1})
 \end{aligned}$$

LSTM in picture

$y^{t+1}$   
 $\uparrow$   
 softmax



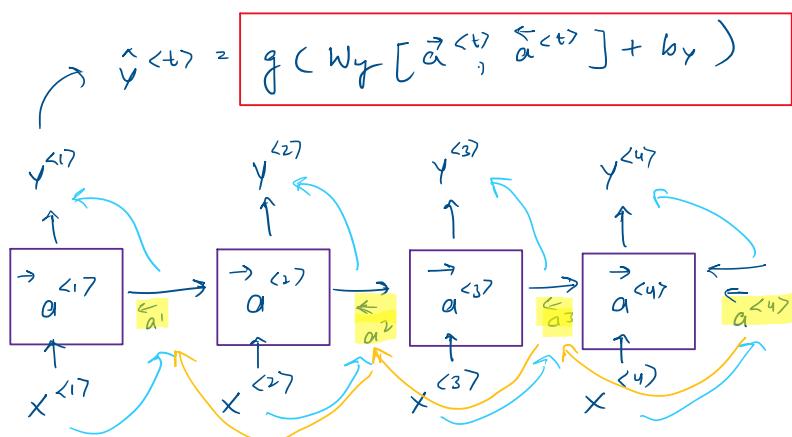
## BIDIRECTIONAL RNN. (BRNN)



He said " Teddy Bears are on sale"

He said "Teddy Roosevelt was a great President"

- Unidirectional RNN.

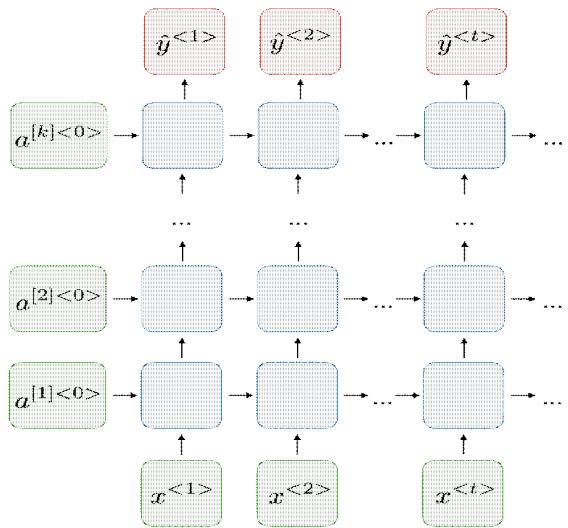


+ word sentence.

Acyclic graph

" BRNN with LSTM "

## DEEP RNN



## MI INTRODUCTION To WORD-EMBEDDINGS.

## WORD REPRESENTATION.

$$V = [q, \text{aaron}, \dots, \text{Zulu}, \text{UNK}]$$

1 hot representation

Man (5391)	Woman (9853)	King (4914)	Queen (7157)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ \vdots \end{bmatrix}$ $D_{(5391)}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ \vdots \\ 1 \\ 1 \end{bmatrix}$ $09853.$	.	- - - -

## Featured Representation

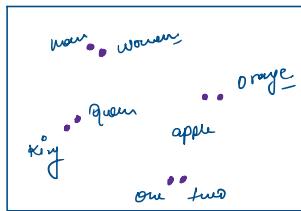
	Man	Women	King	Queen	Apple	Orange
gender	-1	1	-0.95	0.97	0.00	0.01
royalty	-0.01	-0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.61	-0.03	-0.2
foot						
size						
:						
---			--	--	--	--
e591						

" I want a glass of orange juice.

" I want a glass of apple juice.

\* Orange and apple are highly correlated ("fruit").

## Visualizing Word Embeddings

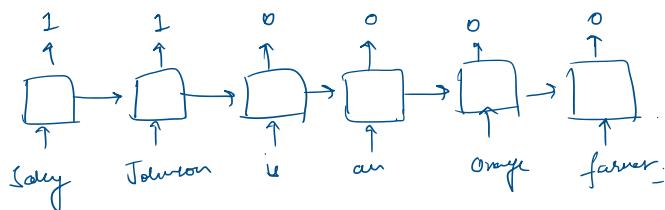


t-SNE.

## USING WORD EMBEDDINGS

" Name Entity Recognition "

Sally Johnson is an Orange farmer.



" Robert is an durian farmer "  
 Robert = a  
 durian = Cultivator.  
 farmer = Cultivator.

"  
 ⇒ Cultivator & Farmer are highly correlated."  
 "

### Transfer Learning

- \* > Learn word embedding from large text-corpus.
- \* > Transfer embeddings to new task with smaller training set.
- > Continue to fine tune the word embeddings with new data.

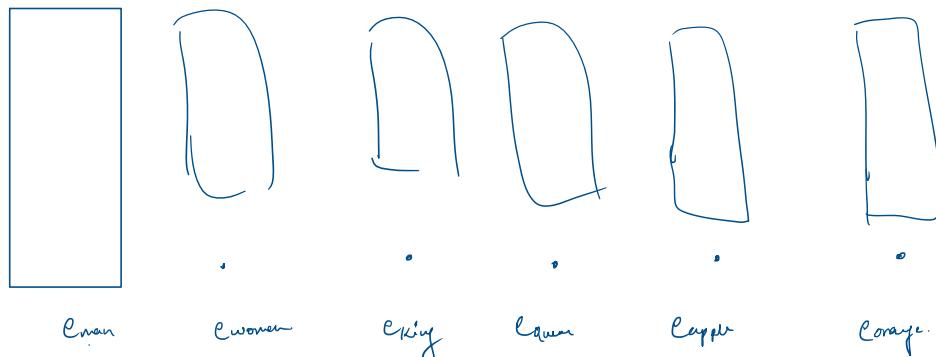
## PROPERTIES OF WORD-EMBEDDINGS

Analogies.

Man : Women :: King : ?

Man	Women	King	Queen	Apple	Orange
(5391)	(9853)	(4914)	(7153)	(456)	(6257)

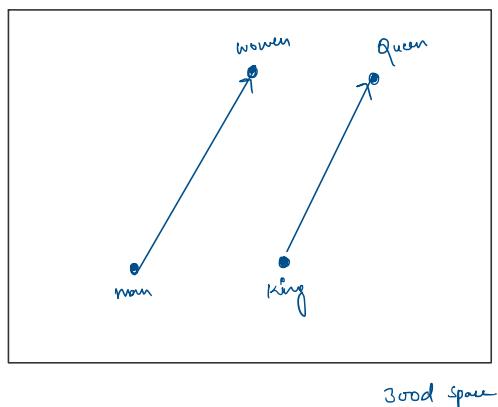
gender  
royalty  
Age  
Food



$$C_{\text{man}} - C_{\text{woman}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad C_{\text{king}} - C_{\text{queen}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$C_{\text{man}} - C_{\text{woman}} \approx C_{\text{king}} - [ \dots ]$$

$$\Rightarrow C_{\text{man}} - C_{\text{woman}} \approx C_{\text{king}} - C_{\text{queen}}$$

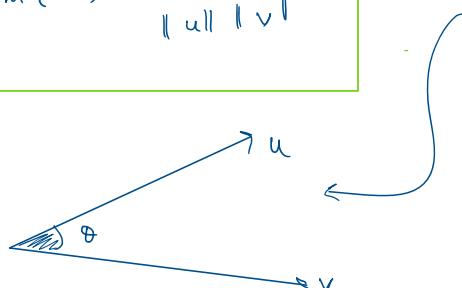


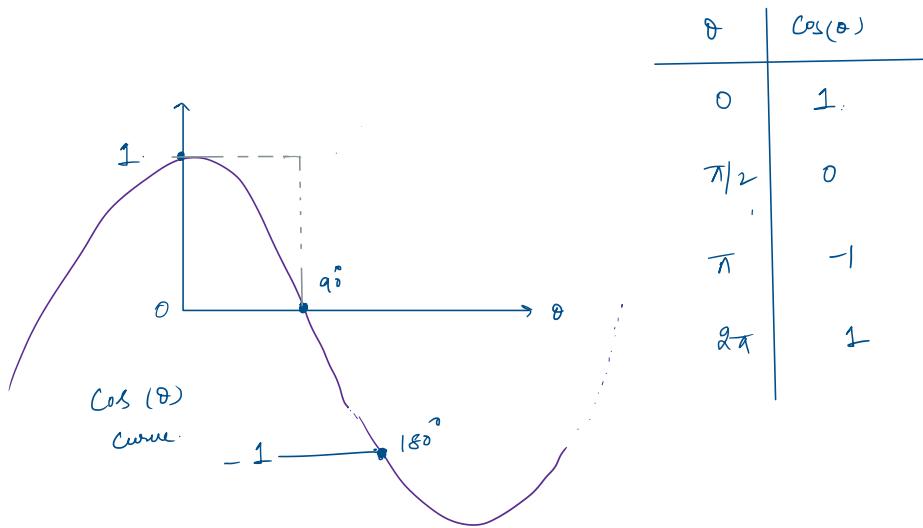
find word  $w$ , s.t:

$$\max \text{Sim.}(ew, C_{\text{king}} - C_{\text{man}} + C_{\text{woman}})$$

- find a similarity fx.  
 $\Rightarrow$  Cosine Similarity.

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|} = \cos(\theta)$$





Other similarity fx:

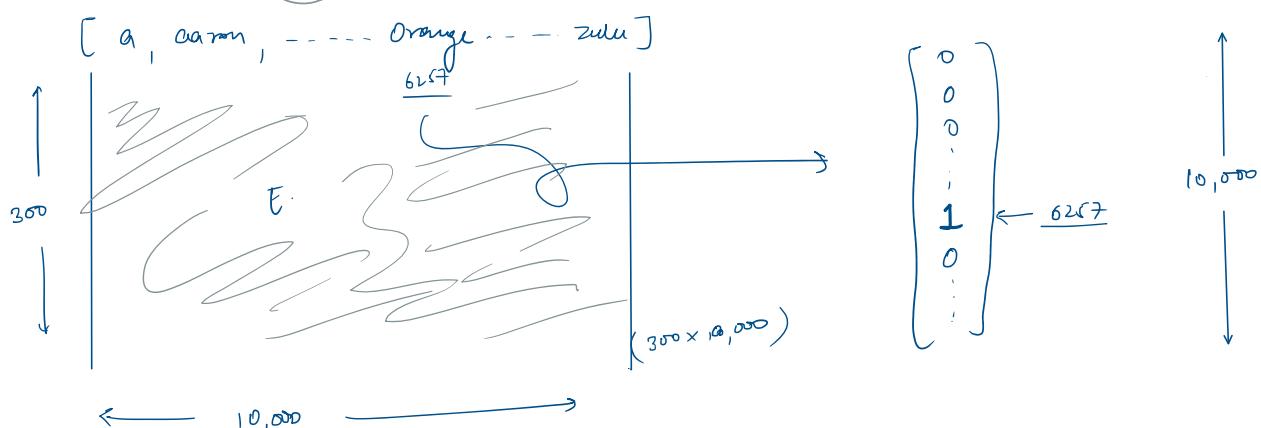
- Euclidean distance.  $\| u - v \|^2$

India : delhi :: Kenya : "Nairobi"

## EMBEDDING MATRIX.

embedding matrix (E)

( $300 \times 10,000$ )  $\rightarrow$  if you have 10,000 word dictionary



$$E \cdot O_{6287} = \begin{bmatrix} \text{apple} \\ \text{orange} \\ \text{carrot} \end{bmatrix}_{300,1} = e_{6287}$$

one-hot vector.

$$\boxed{E \cdot O_j = e_j}$$

Embedding for word  $j$ .

Embedding for word.

\* It is not efficient to use multiplication of matrix, as the cost is very high.

## MD LEARNING WORD EMBEDDINGS : Word2Vec and Glove

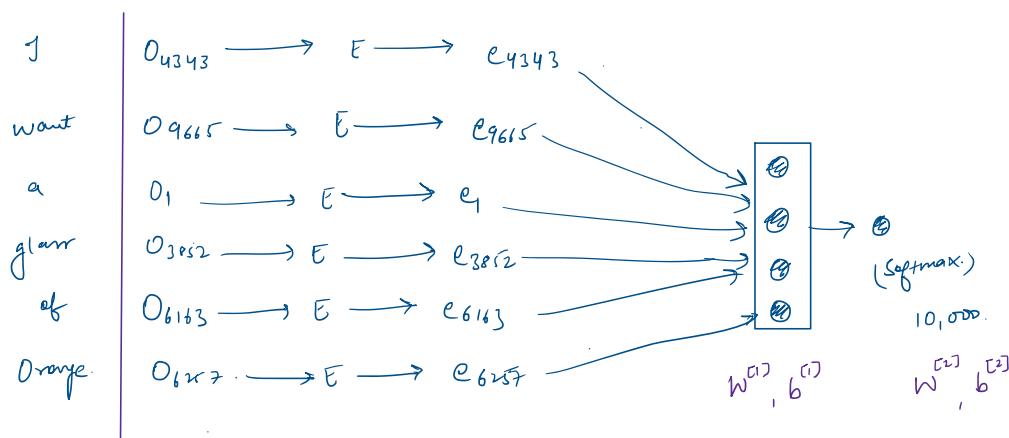
### LEARNING WORD EMBEDDING

Natural Language Model.

" I want a glass of Orange \_\_\_\_\_ .

(4343) (9665) (3852) (6163) (6257)

.....  
index in vocab.



### WORD2VEC

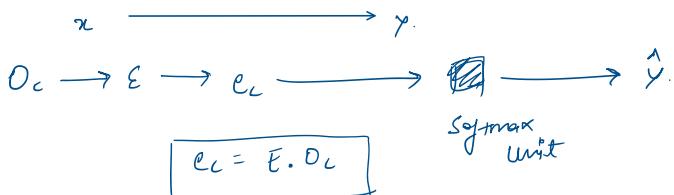
Skip-grams : I want a glass of orange juice along with my cereal.

Context	Target
randomly	juice
selected	glass
	my

Model

$$\text{vocab-size} = 10,000 \text{K}$$

Context  $c$  ("Orange")  $\longrightarrow$  Target  $t$  ("juice")  
 $(6257)$   $(4834)$



Softmax :  $p(t|c) = \frac{e^{\theta_t^T \cdot e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T \cdot e_c}}$

$\theta_t$  = parameter associated with O/P  $t$   
 (class of  $t$ , being a label).

$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{1000} [y_i \cdot \log \hat{y}_i]$$

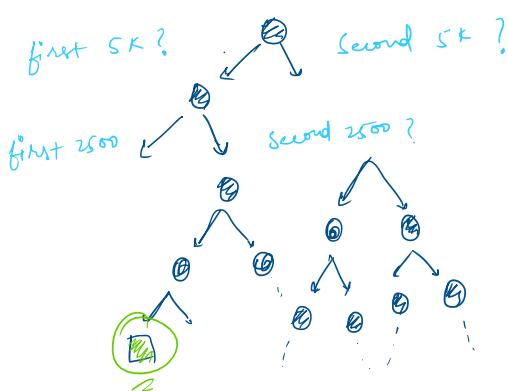
### Problems with Softmax Classification.

- Computational Speed.

$$p(t|c) = \frac{e^{\theta_t^T \cdot e_c}}{\sum_{i=1}^{10,000} e^{\theta_i^T \cdot e_c}}$$

every time it will calculate 10k times,  
 at higher values, it will compute  
 slower.

$\Rightarrow$  Hierarchical softmax.



### NEGATIVE SAMPLING.

Predict (Context, target) pair

Context	target	<u>Probability</u>
Orange	juice	$\approx 1$
Orange	king	$\approx 0$
Orange	book	$\approx 0$

$K$

*find p with random K target and fill them (0).*

> Create a supervised learning algorithm, given (Context, target) predict probability.

$$K = \begin{cases} 5-20, & \text{for small datasets} \\ 2-5, & \text{for larger datasets} \end{cases}$$

Softmax Model  $\rightarrow$

$$p(t|c) = \frac{e^{\theta_t^T \cdot e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T \cdot e_c}}$$

Context	word	target
- - -	- - -	1
- - -	- - -	0
- - -	- - -	0
- - -	- -	0

use Logistic Regression Model.

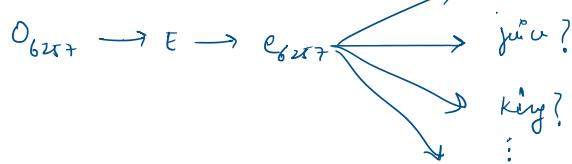
$$p(y=1 | c, t) = \sigma(\theta_t^T \cdot e_c)$$

Sigmoid

"Orange"

6257

$e_{6257}$



1, real ex.  
k, random ex.

Selecting -ve examples?

$$P(w_i) = \frac{f(w_i)^{1/4}}{\sum_{j=1}^{10,000} f(w_j)^{1/4}} \approx \frac{1}{|V|}$$

uniform distribution raised to  $1/4^{th}$ .

$$\sum_{i=1}^{10,000} f(w_j)^{\frac{1}{k}}$$

|V|

to  $S^T$

## GLOVE WORD VECTORS

global vectors for word representation.

"I want a glass of Orange juice to go along with my cereal."

$x_{ij} = \# \text{ of times } j \text{ appears in the context of } i$

$$\text{minimize } \sum_{i=1}^{10k} \sum_{j=1}^{10k} f(x_{ij}) \left( \theta_i^T e_j + b_i + b_j - \log x_{ij} \right)^2$$

$\boxed{i, j = t, c}$

or  $(\theta_t)^T e_c$

Correlation

$f(x_{ij})$   
 this is the of at .....  
 (stop words)

\* might give more weight to  
 stop words and lure on more attention. words.

if  $(x_{ij} = 0)$  then  $\log 0 = \text{undefined}$ , so we will add  
 weighing term  $f(x_{ij})$ ,  $f(x_{ij}) = 0$ , if  $\boxed{x_{ij} = 0}$ .  
 $\therefore 0 \cdot \log 0 = 0$ .

$f$  can be chosen to be heuristic to accomplish this.

Take average after calculating.

$$e_w^{(\text{final})} = \frac{e_w + \theta_w}{2}$$

## APPLICATIONS USING WORD EMBEDDINGS.

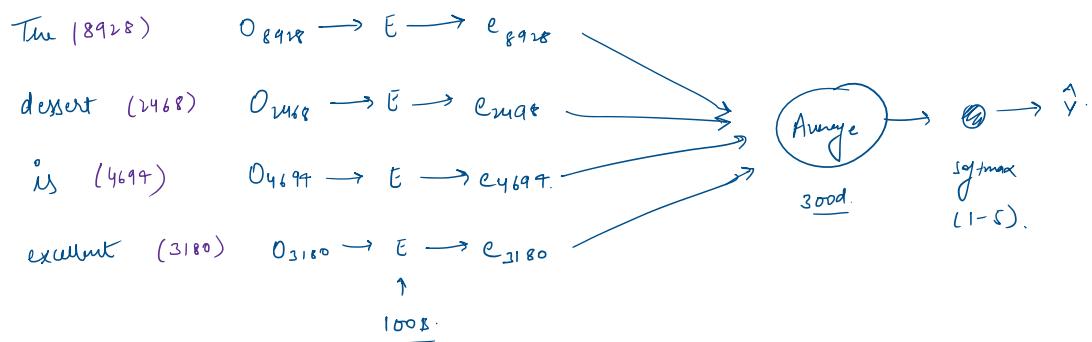
### SENTIMENT CLASSIFICATION.

$$x \longrightarrow y$$

"dumb is excellent"

5 \*

" dessert is excellent"	5*
" Service is slow"	1*
" good for quick meal, but nothing special"	3*



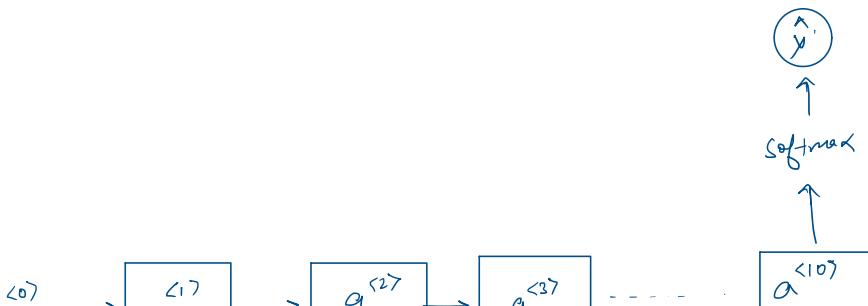
- > ignores word order.
- = " Completely lacking in good taste, good service and good ambience."
- \* negative review, but "good" appears a lot!

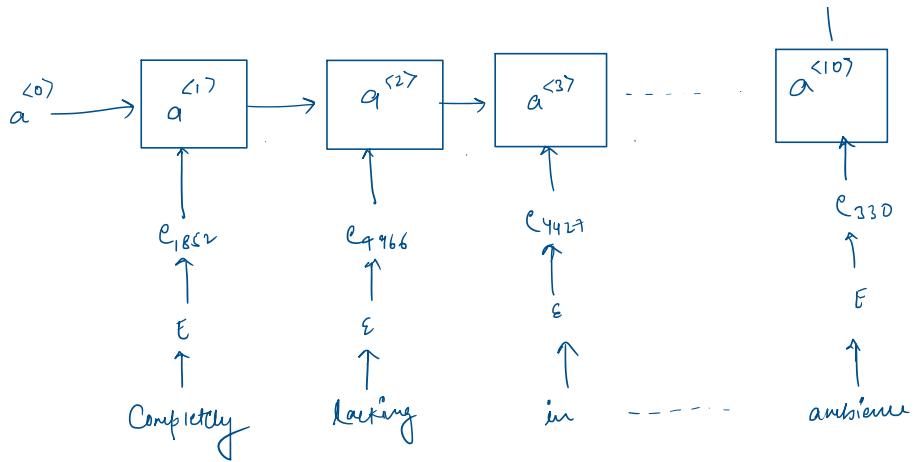
## # RNN for Sentiment Classification.

Completely lacking - - - - - - - - - ambience.



- ① Calculate one-hot vector.
- ② Embedding Matrix
- ③ find embedding vector. ( $O \times \bar{E}$ )
- ④ feed  $e$  to RNN





"Many to One?"

## DEBIASING WORD EMBEDDINGS

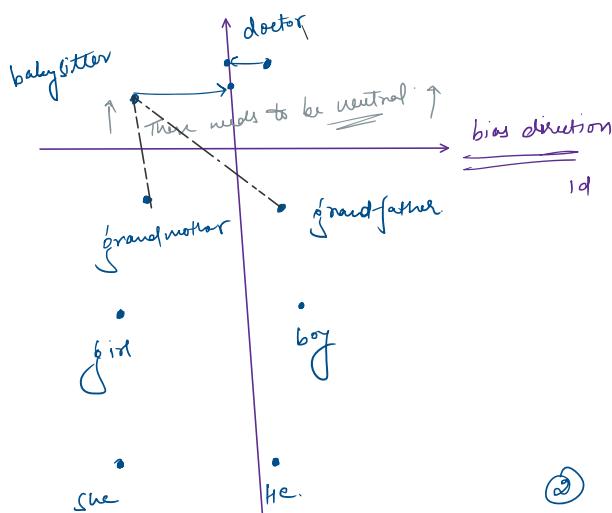
Man : Women :: King : Queen.

Man : Programmer :: Women : Homemaker.  $\times$

Father : Doctor :: Mother : Nurse  $\times$

\* Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of text used to train the model.

non-bias direction.  $\overrightarrow{w_{id}}$



① Identify bias direction

(Taking gender bias here.)

$$\begin{aligned} & \text{E}_\text{men} - \text{E}_\text{women} \\ & \text{E}_\text{male} - \text{E}_\text{female} \\ & \text{E}_\text{boy} - \text{E}_\text{girl} \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Average}$$

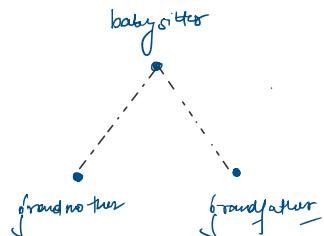
② Neutralize: For every word that is not definitional, project to get rid of bias.

③ Equalize pairs :-



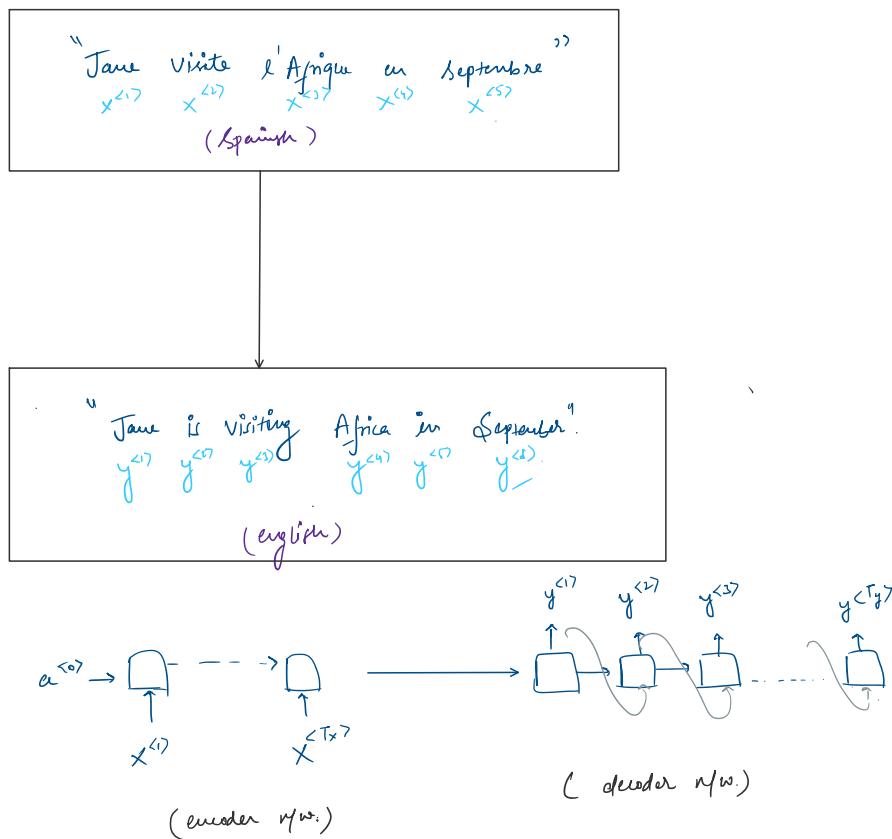


\* Move grandmother and grandfather at pair points equidistant from a axis

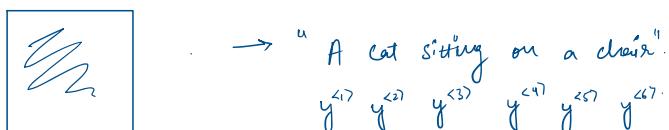


## ML: VARIOUS SEQUENCE TO SEQUENCE ARCHITECTURES.

### Basic Model

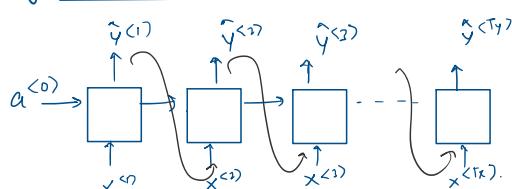


### Image Captioning

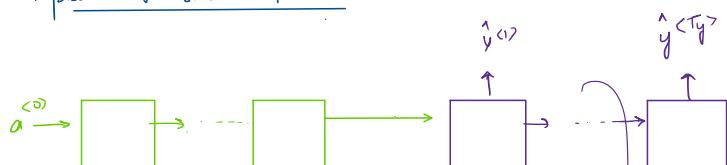


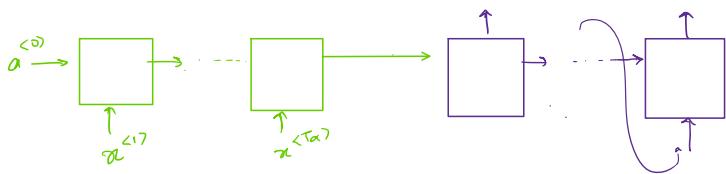
### Picking the most likely sentence

Language Model (as in first week)



### Machine Translation Model





encoded n/w.

decoded n/w.

Conditional N/W.

"Jane visited Africa in September"

$$P(y_1, y_2, y_3, y_4, \dots)$$

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September Jane will visit Africa.
- Her African friend welcomed Jane in September.

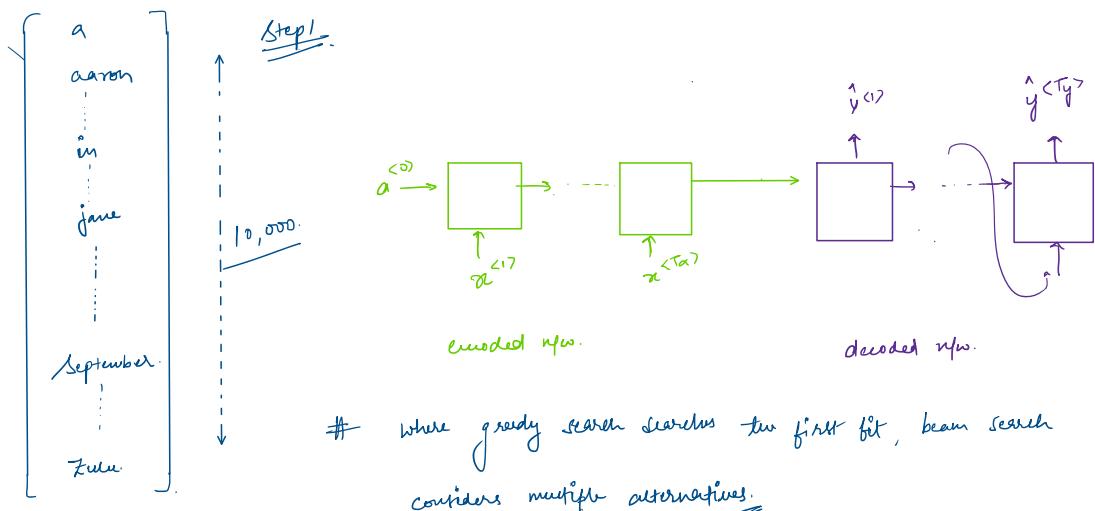
good translation

bad translation

"beam search"

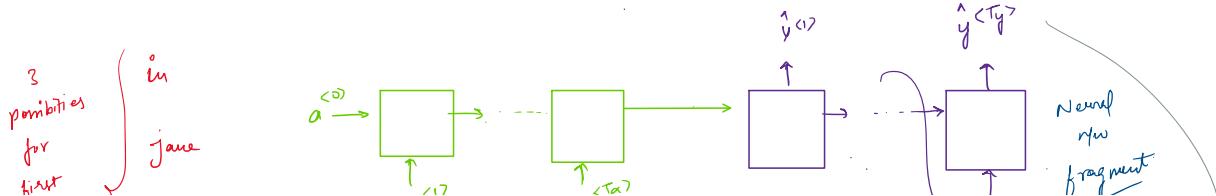
## Beam Search.

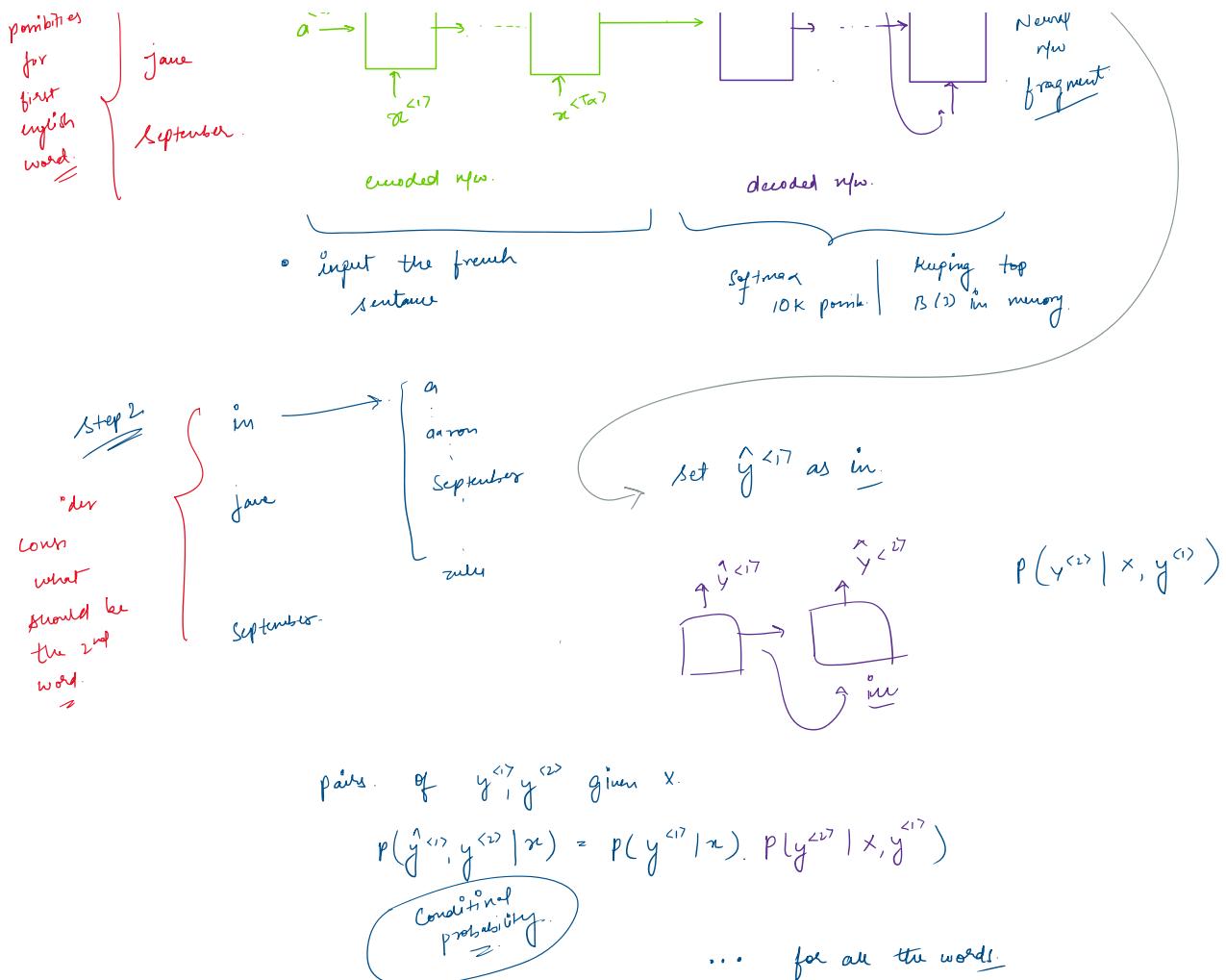
### > Beam Search Algorithm



- \* beam search algorithm has a perimeter  $B$  (Beam width.)

$B=3$  (say).    • Considers 3 possibilities.





## Refinements to Beam Search

### Length Normalization

$$\text{argmax}_y \prod_{t=1}^{T_y} p(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

$\log$

$$p(y^{<1>} \dots y^{<T_y>} | x) = p(y^{<1>} | x) \cdot p(y^{<2>} | x, y^{<1>}) \cdot p(y^{<3>} | x, y^{<1>}, y^{<2>})$$

\* Normalize by number of words.

$$\frac{1}{T_y} \sum_{t=1}^{T_y} \log p(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

Beam Width?

Larger :- Higher Accuracy, more Computational

Smaller :- Low Accuracy , Faster

## Error Analysis On Beam Search.

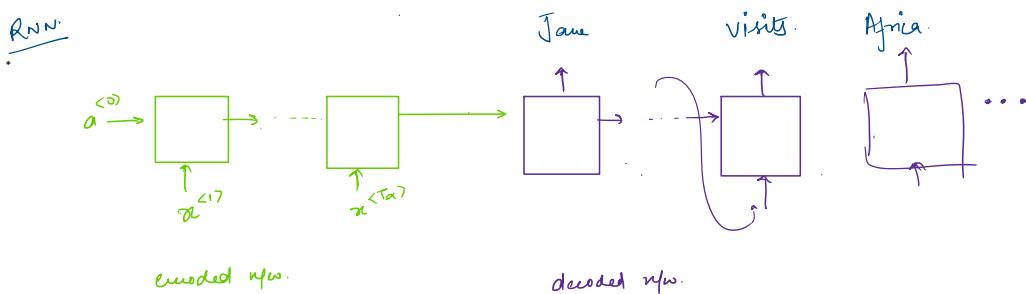
- # Beam search is an approximate search algorithm. / heuristic search algorithm.

<sup>11</sup> Jane visite l'Afrique en septembre. <sup>12</sup>

Theme : Jane visits Africa in September. (y\*)

Algorithm : Jane visited Africa last September. (y)

- RNN (encoder and decoder)
  - Beam Search.



$$\text{mod. } \left( p(y^* | x), \quad p(y | x) \right)$$

Jane visits Africa in September. (y\*)

$$P(y^* | x)$$

Jane visited Africa last September. (iy)

$$P(\hat{y} | x)$$

Case 1 :  $P(y^*) > P(\hat{y})$

$\Rightarrow$  but beam search chooses  $P(\hat{y})$ .

⇒ Beam Search is at fault

Case 2 :  $P(y^*) < P(\hat{y})$

$\Rightarrow$   $y^*$  is a better translation than  $\hat{y}$ . But RNN predicts

$$p(y^*) < p(\hat{y})$$

⇒ RNN model is at fault.

## Bleau Score

( Bilingual Evaluation Understudy Score. )

- given a french sentence there can be multiple english translation,  
this can be solved with "Bleau Score"

$\alpha - \alpha - \alpha - \alpha$

French : Le chat est sur le tapis.

Refrene1 : The cat is on the mat.

Refrene2 : There is cat on the mat.

BLEAU Score measure how good a machine translation is.

$$\text{Precision} := \frac{2}{7}$$

## Bleau Score on Bigrams

"pairs of words next to each other."

Refrene1 : The cat is on the mat.

Refrene2 : There is cat on the mat.

MT O/P : The cat the cat on the mat

(Machine Translation)

	MT O/P	Refrene1
the cat	Count	2
Cat the	Count	1
Cat on	Count	1
on the	Count	1
the mat.	Count	1
	CountClip	1 0 1 1 1

Modified bigram Precision :-

$$\frac{\sum \text{Count Clips}}{\sum \text{Count}}$$

$$= \frac{4}{6} = \frac{2}{3}$$

### Bleu Details

$P_n$  = Bleu Score of  $n$ -grams.

$P_1, P_2, P_3, P_4$

Combined bleu score =  $BP \exp\left(\frac{1}{4} \sum_{n=1}^4 P_n\right)$ .

(brevity penalty)

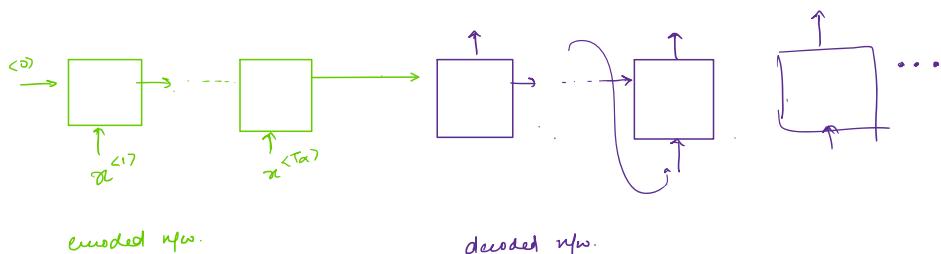
$$BP \begin{cases} 1, & \text{if MT-Output length} > \text{reference O/P Length} \\ \exp\left(1 - \frac{\text{reference-O/P Length}}{\text{MT-O/P Length}}\right) & \end{cases}$$

### Attention Model Intuition

#### # The problem of long Sequence

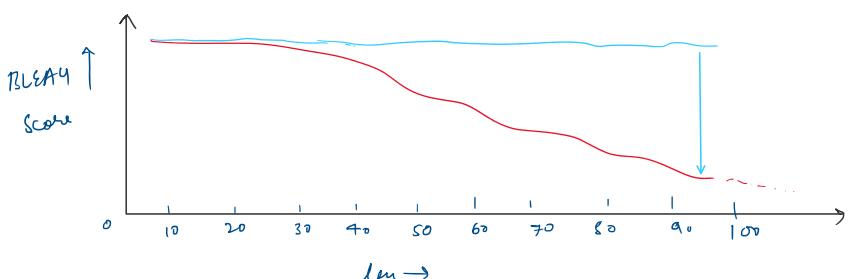
> given a long french sentence

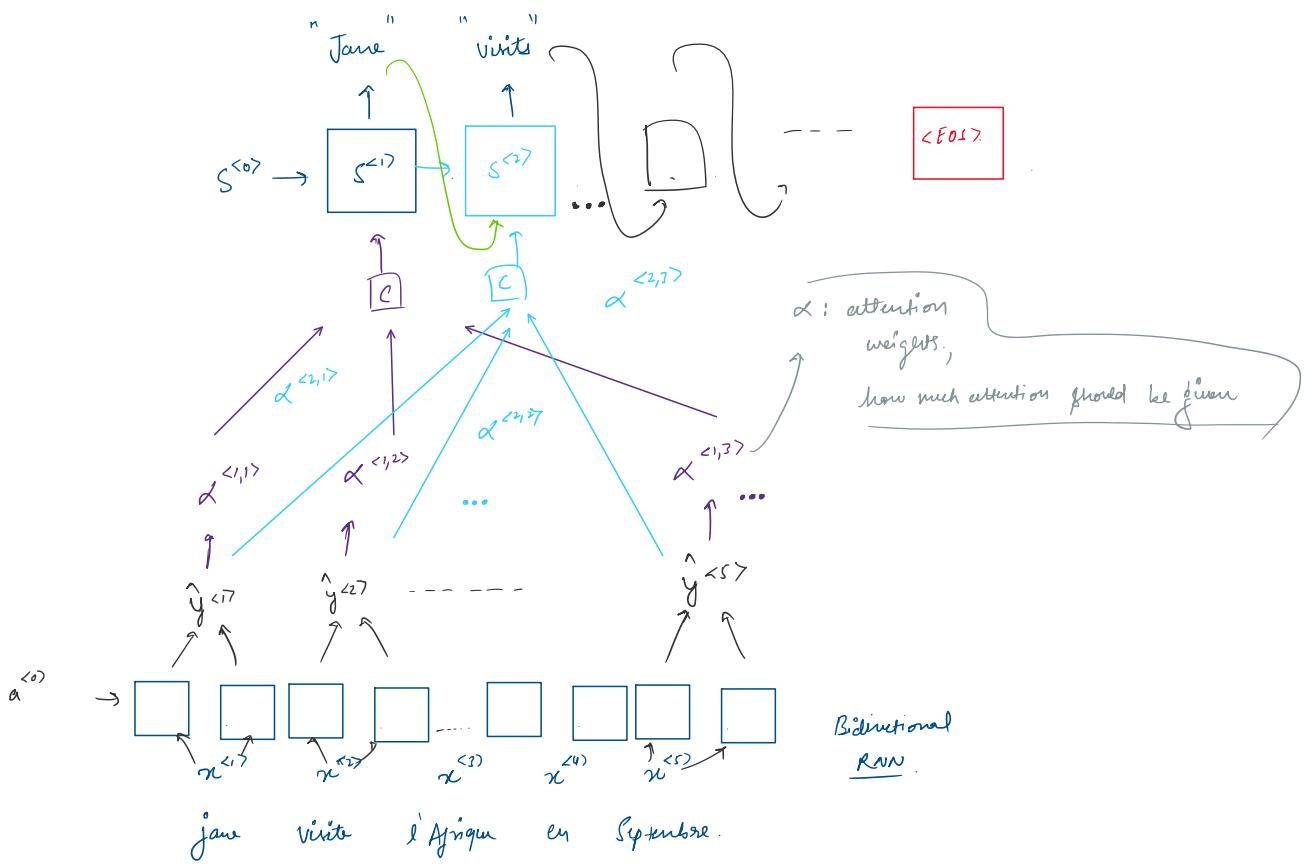
-----  
-----



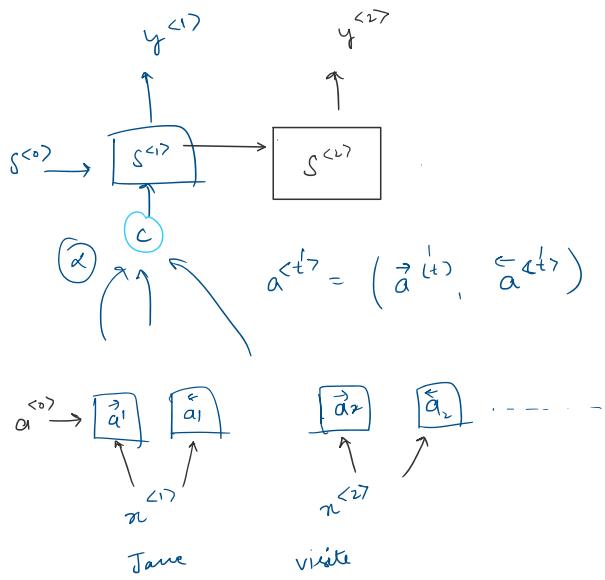
\* It is hard to memorize the whole sentence, what neurons will do is

translate a part of sentence and move forward.





## Attention Model



$$\sum_{x^1} \alpha^{<1, x^1>} = 1$$

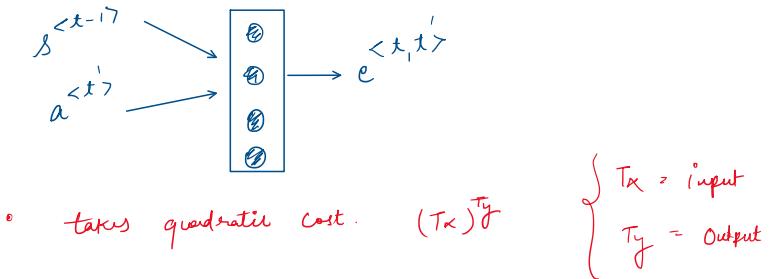
$\alpha^{<t, t>}$  = amount of attention  $y^{<t>}$  should pay  
to  $a^{<t>}$

## Competing Attention

$a^{<t,t'>}$  = amount of attention  $y^{(t)}$  should pay to  $a^{<t'>}$

$$\langle t, \lambda' \rangle \quad \exp \left( e^{\langle t, \lambda' \rangle} \right)$$

$$\alpha^{(t, t')} = \frac{\exp(e^{(t, t')})}{\sum_{t'=1}^T \exp(e^{(t, t')})}$$



$$\left\{ \begin{array}{l} T_x = \text{input} \\ T_y = \text{output} \end{array} \right.$$

## M2 Speech Recognition

### Speech recognition

$$x \longrightarrow y$$

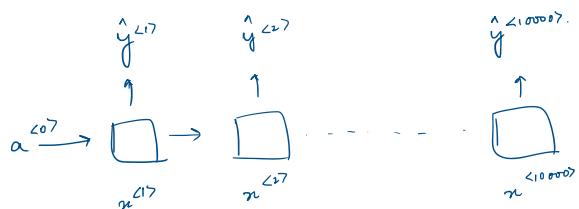
(audio clip)

(transcript)



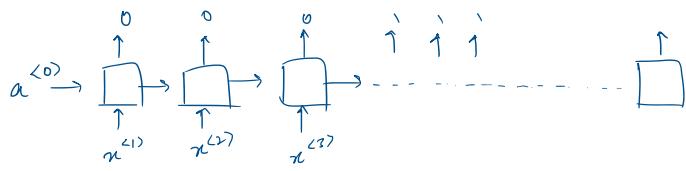
### CTC Cost

### "Connectionist temporal Classification"



### Trigger Word detection

Spectrogram



"Hey Siri !!"

"OK google !!"

[...].